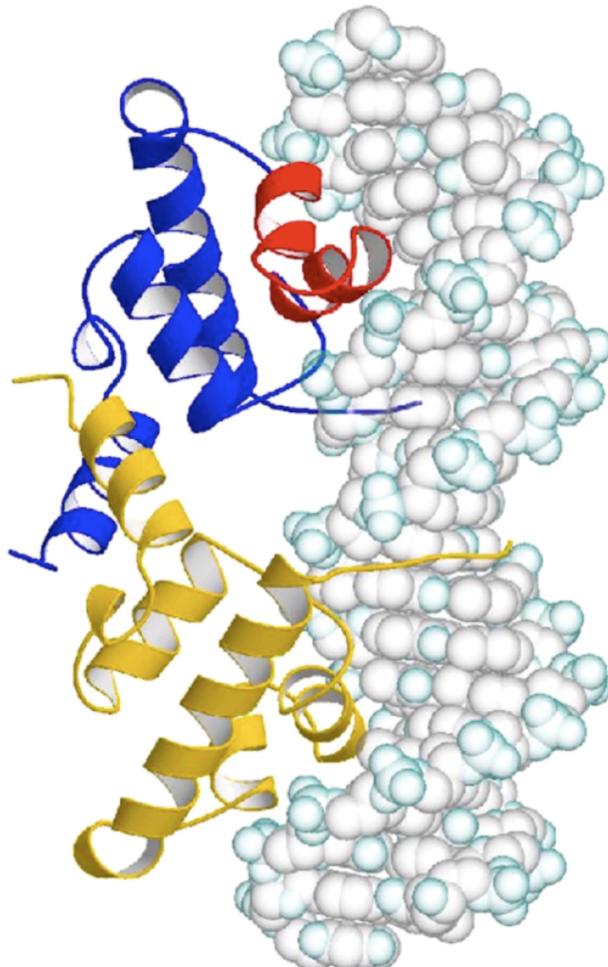




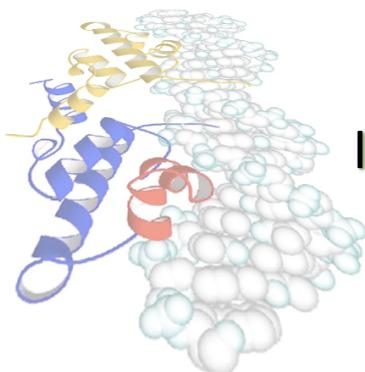
MAHIDOL
UNIVERSITY
Wisdom of the Land



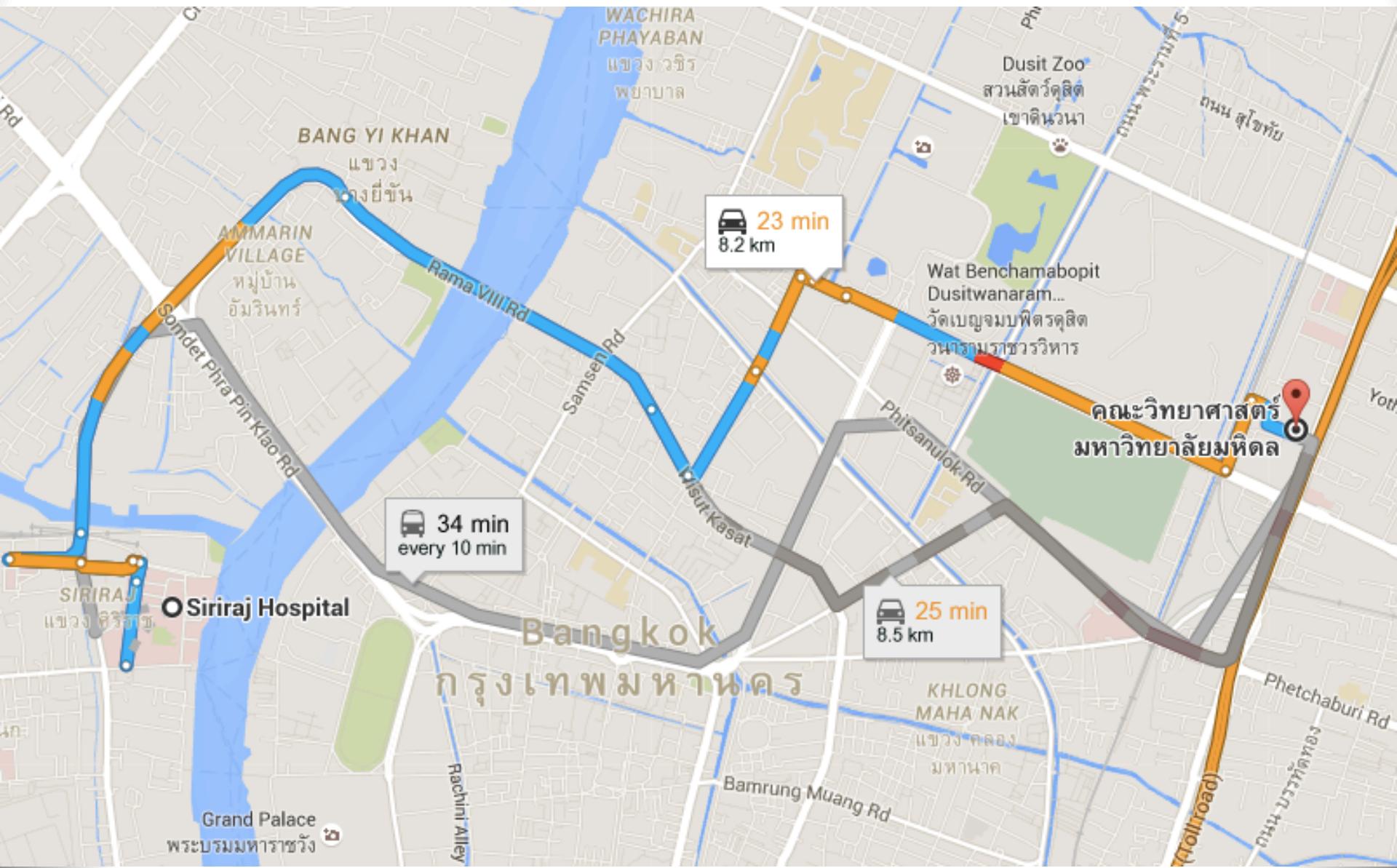
Systems Biology of Transcriptional Regulation

Varodom Charoensawan

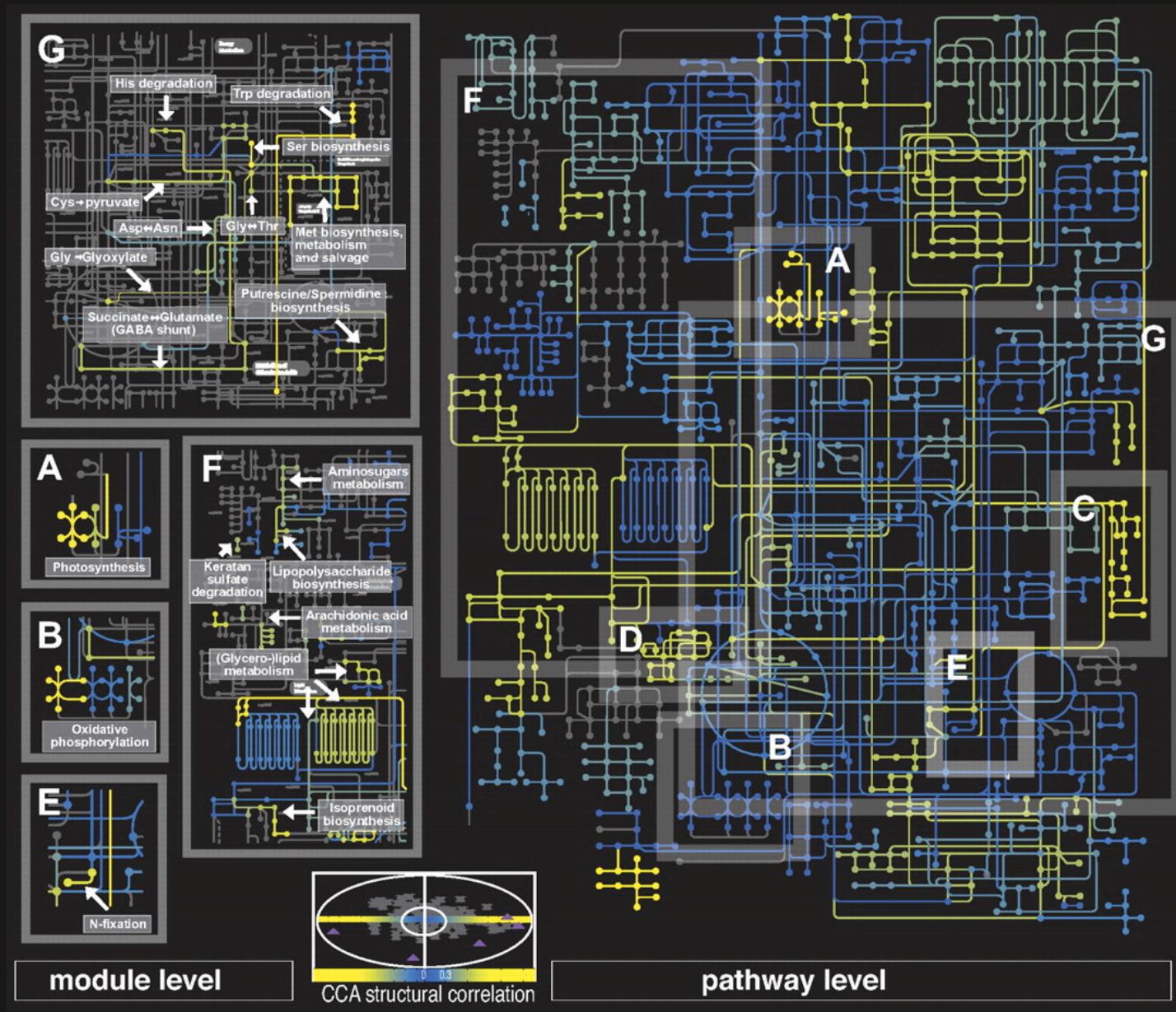
Department of Biochemistry, Faculty of Science
Integrative Computational Bioscience (ICBS) Center



How to get from A to B?



A complete map of an organism?



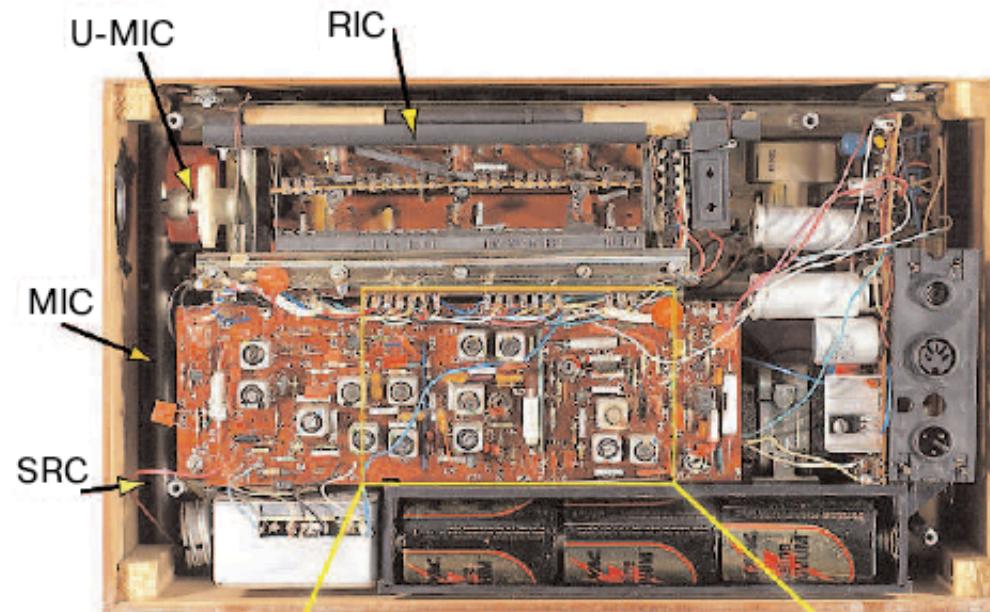
Gianoulis et al. PNAS (2008) 106(5): 1374–1379,

“Systems Biology”

- Investigates links between multiple genes, proteins at one time, rather than one individual gene, protein (genome-wide, large-scale and –omics data)
- Usually involved combination of different techniques, tools: wet-lab experiment + computational analyses

Creating a biological “map” of an organism

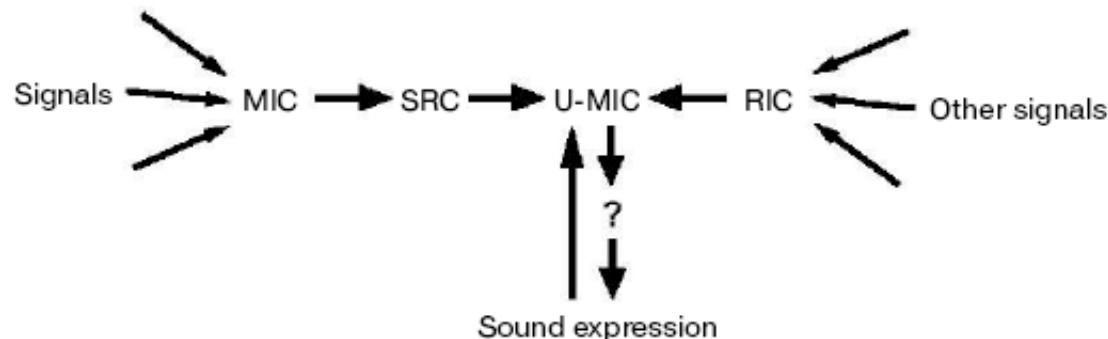
Can a biologist fix a radio? (1)



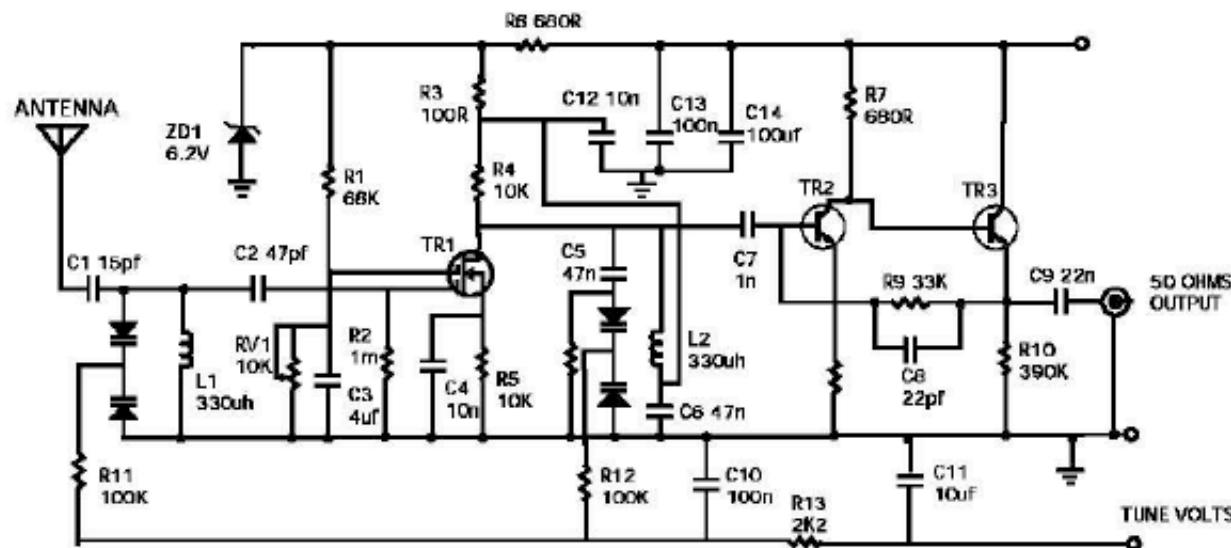
- Serendipitously Recovered Component
- Really Important Component
- Most Important Component
- Undoubtedly Most Important Component

Can a biologist fix a radio? (2)

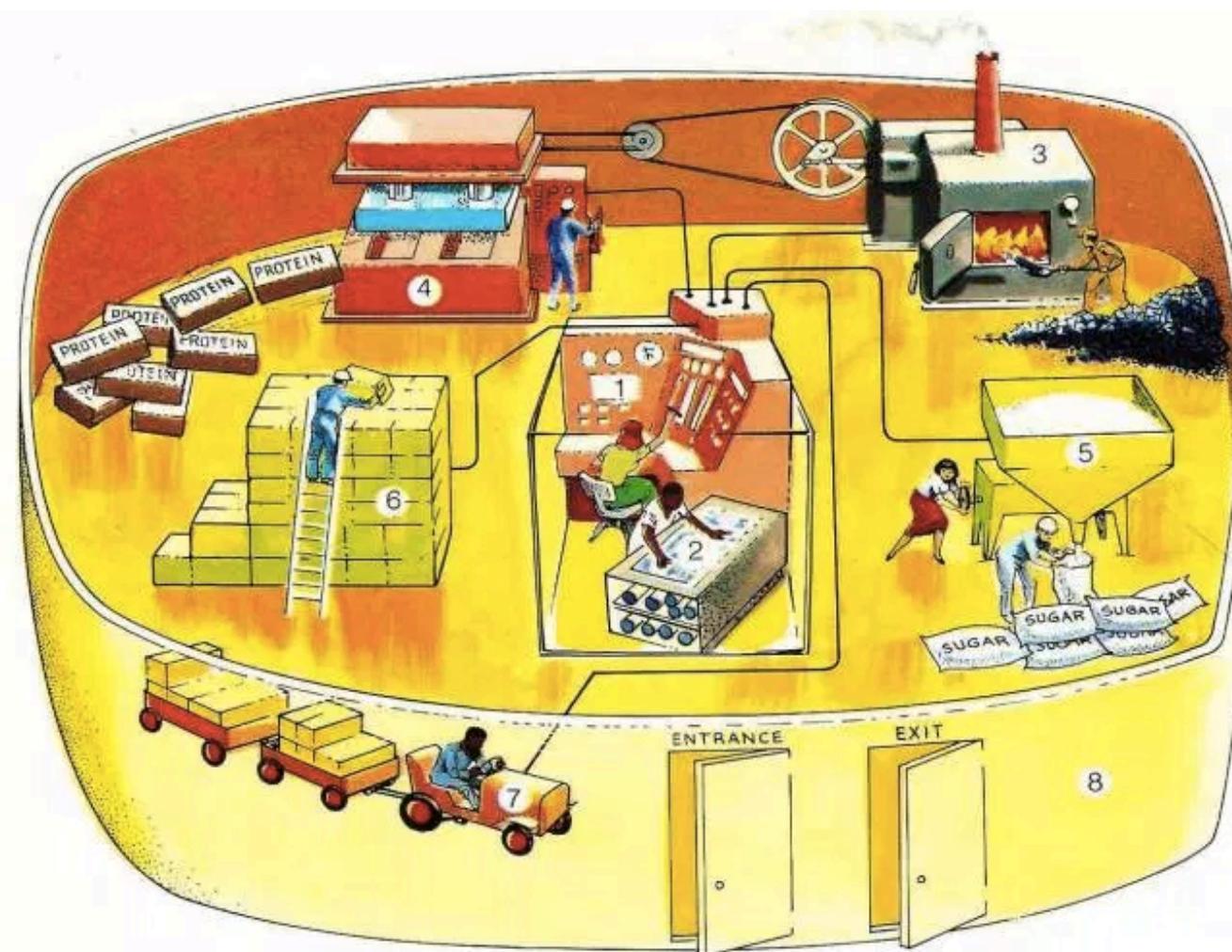
a



b

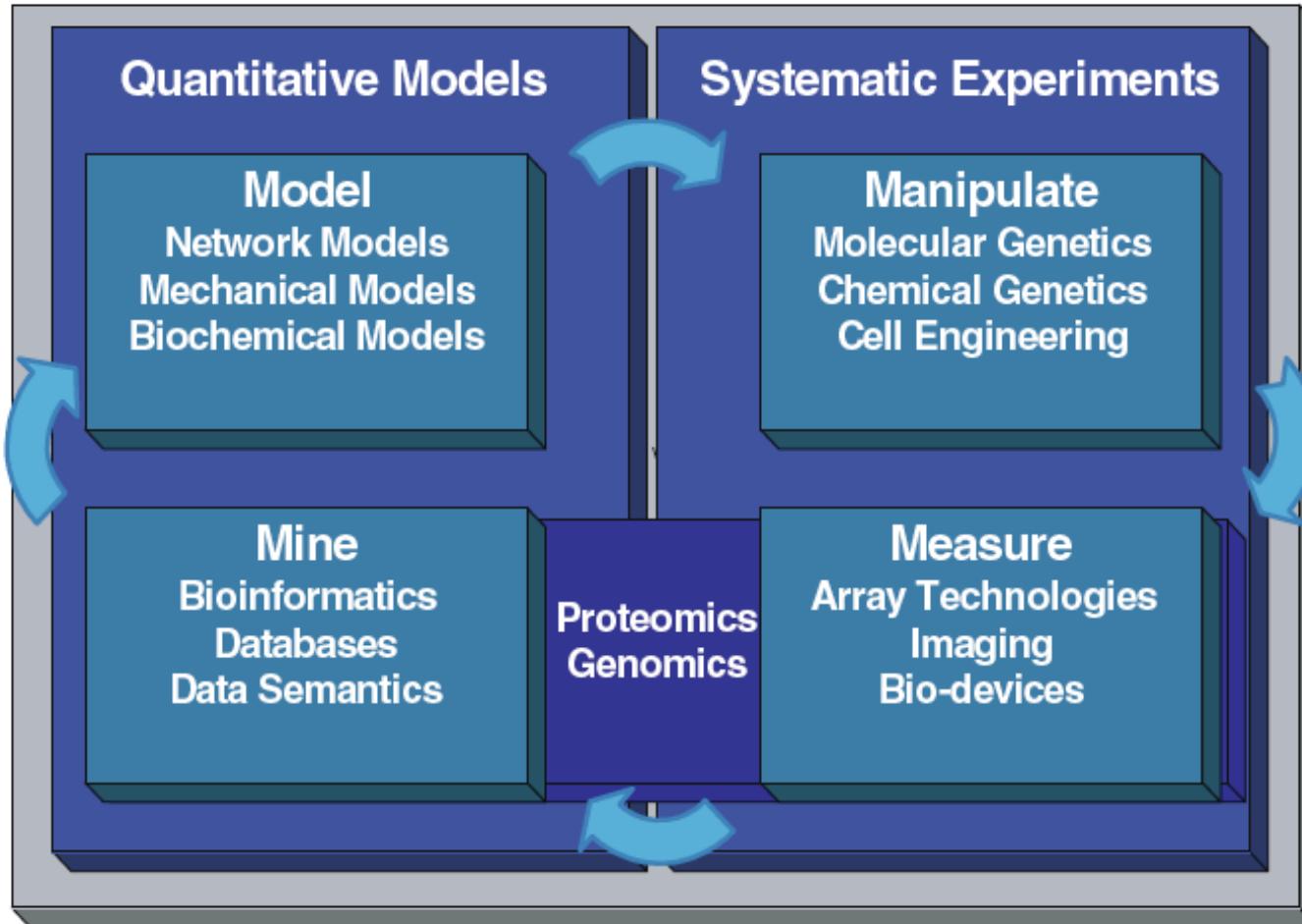


Cell Factory



1. Nucleus (control center)
2. Chromosomes (file cabinet, contains all the information)
3. Mitochondria (power house)
4. Ribosomes and RER (assembly line, work benches)
5. Chloroplast (energy source)
6. Golgi Apparatus (packaging center)
7. Vesicles (transportation of products)
8. Cell Membrane (security gate)

4Ms: operational definition



Typical Workflow of Bioinformatics

The screenshot shows the Ensembl genome browser interface for Human (GRCh37). The top navigation bar includes links for BLAST/BLAT, BioMart, Tools, Downloads, Help & Documentation, Blog, and Mirrors. The user is currently viewing the location 13:32,889,611-32,973,347, specifically the Gene: BRCA2 Transcript: BRCA2-001. The left sidebar contains a tree menu for "Location-based displays" with categories like Whole genome, Chromosome summary, Region overview, Region in detail, Comparative Genomics, Genetic Variation, and Other genome browsers. Below this are buttons for Configure this page, Manage your data, Export data, and Bookmark this page. The main content area displays the "Region in detail" view for Chromosome 13. It features a genomic track with bands representing different genomic features. A red box highlights a specific region on the q13.3 band. The track shows various genes and transcripts, including EEF1DP3, FRY, RP11-207N4.3, ZAR1L, BRCA2, RP11-37E23.5, IFT1P1, ATP8A2P2, N4BP2L1, N4BP2L2, and SNORA16. The forward strand is indicated by a blue arrow at the top right. A legend at the bottom defines the colors: blue for processed transcript, grey for pseudogene, and purple for RNA gene. An "Export Image" button is located at the bottom right of the main panel.

- Organize / store / retrieve biological information from databases and public domains
- Perform some basic statistical testing

“Synthetic Biology”

- A discipline that uses engineering principles to design and assemble improved biological components and systems
- Artificially design and construction of biological devices systems, and machines
- Combining interdisciplinary methods from biotechnology, genetic engineering, molecular biology, biophysics

Biobricks

- Standardized biological parts are DNA sequences
- Common interface and are designed to be composed and incorporated into living cells

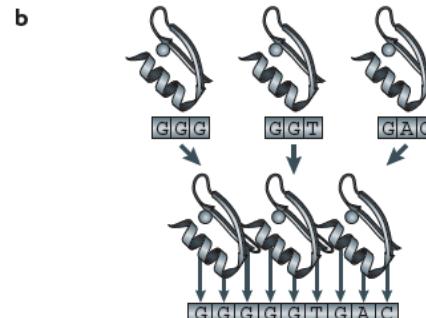
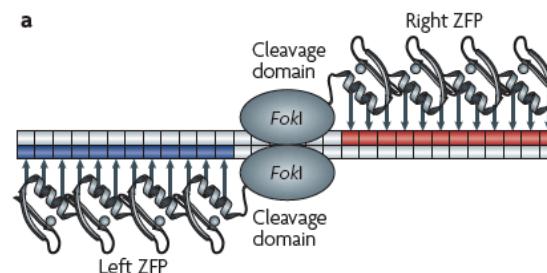


Biobricks

	promoter		primer binding site
	cds		restriction site
	ribosome entry site][blunt restriction site
	terminator	L	5' sticky restriction site
	operator	L	3' sticky restriction site
	insulator	—	5' overhang
	ribonuclease site	—	3' overhang
	rna stability element	—	assembly scar
	protease site	x—	signature
	protein stability element	——	user defined
	origin of replication		

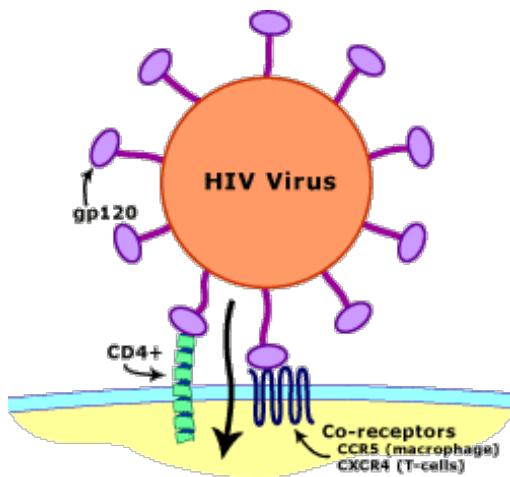
Genome Editing and Reverse Engineering

- Zinc finger nucleases (ZFNs)
- Transcription Activator-Like Effector Nucleases (TALENs)
- CRISPR/Cas system



Establishment of HIV-1 resistance in CD4⁺ T cells by genome editing using zinc-finger nucleases

Elena E Perez^{1,2}, Jianbin Wang³, Jeffrey C Miller³, Yann Jouvenot^{3,4}, Kenneth A Kim³, Olga Liu¹, Nathaniel Wang³, Gary Lee³, Victor V Bartsevich³, Ya-Li Lee³, Dmitry Y Guschin³, Igor Rupniewski³, Adam J Waite³, Carmine Carpenito¹, Richard G Carroll¹, Jordan S Orange², Fyodor D Urnov³, Edward J Rebar³, Dale Ando³, Philip D Gregory³, James L Riley¹, Michael C Holmes³ & Carl H June¹



Pict from: walmartramen.blogspot.com/2014/03/bone-marrow-transplant-and-engineered.html

CRISPR in 2018: Coming to a Human Near You

The first clinical trials are slated to begin in the U.S. and Europe while others are stalled.

by Emily Mullin December 18, 2017

<https://www.technologyreview.com/s/609722/crispr-in-2018-coming-to-a-human-near-you/>



BETTER FASTER CRISPR

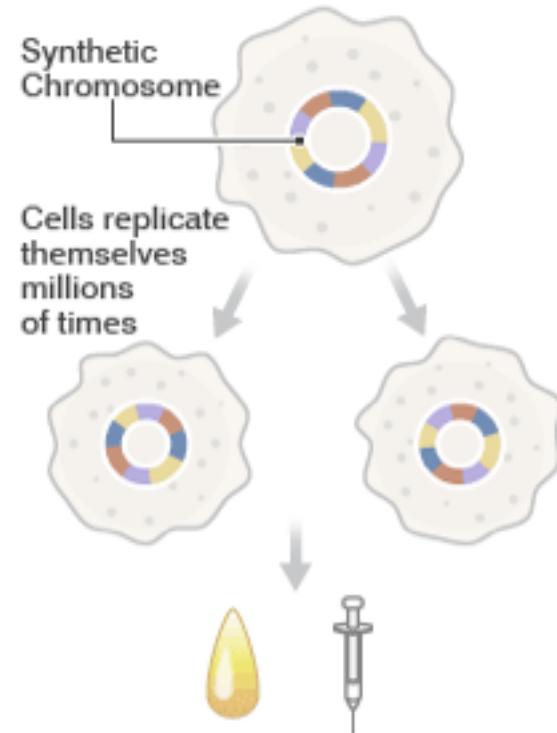
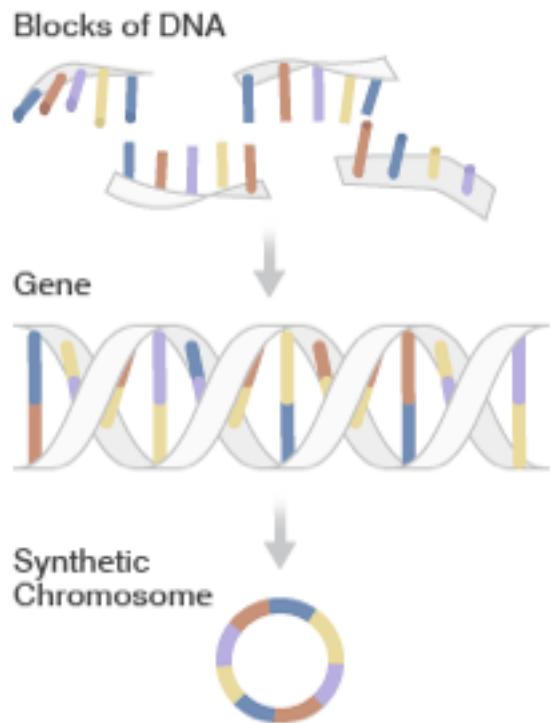
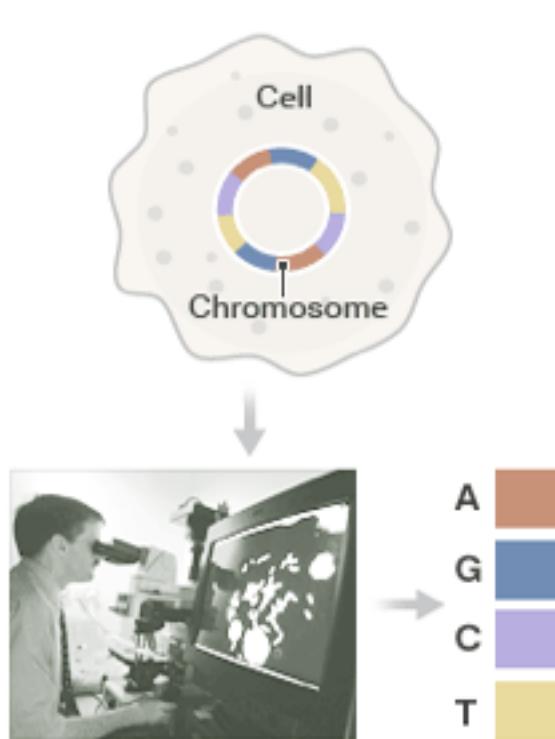
Chinese scientists used Crispr gene editing on 86 human patients

By Katherine Ellen Foley · January 23, 2018

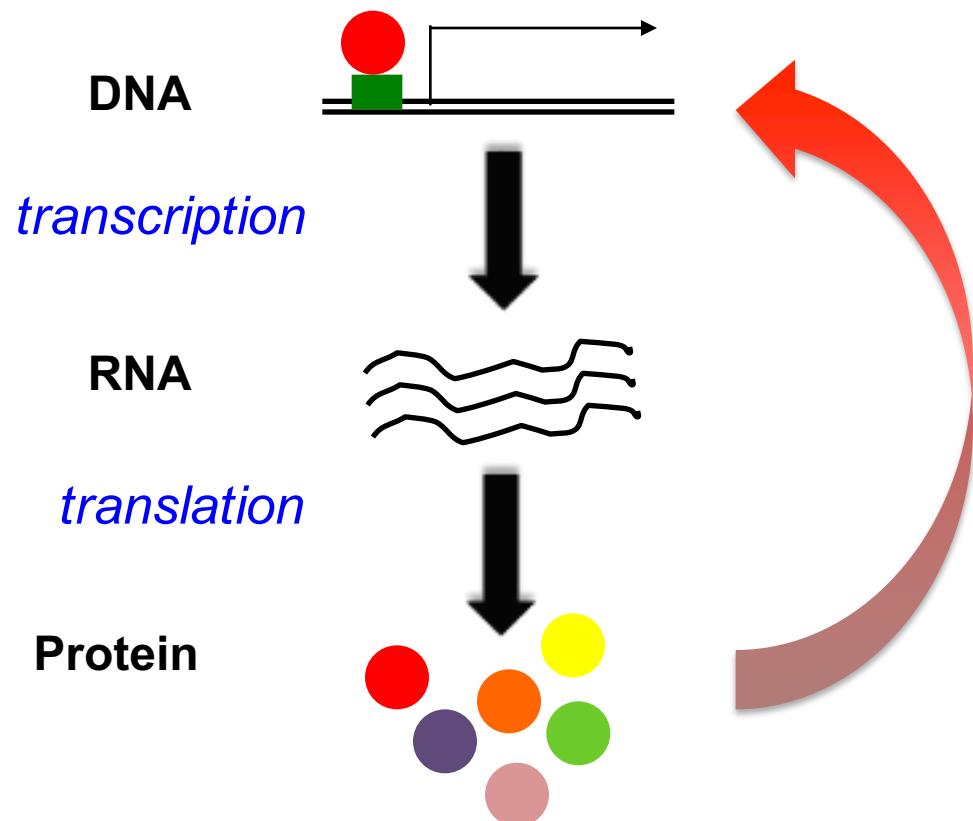


<https://qz.com/1185488/chinese-scientists-used-crispr-gene-editing-on-86-human-patients/>

Venter's artificial life?



The world of “Omics” data



Genome size (bp)

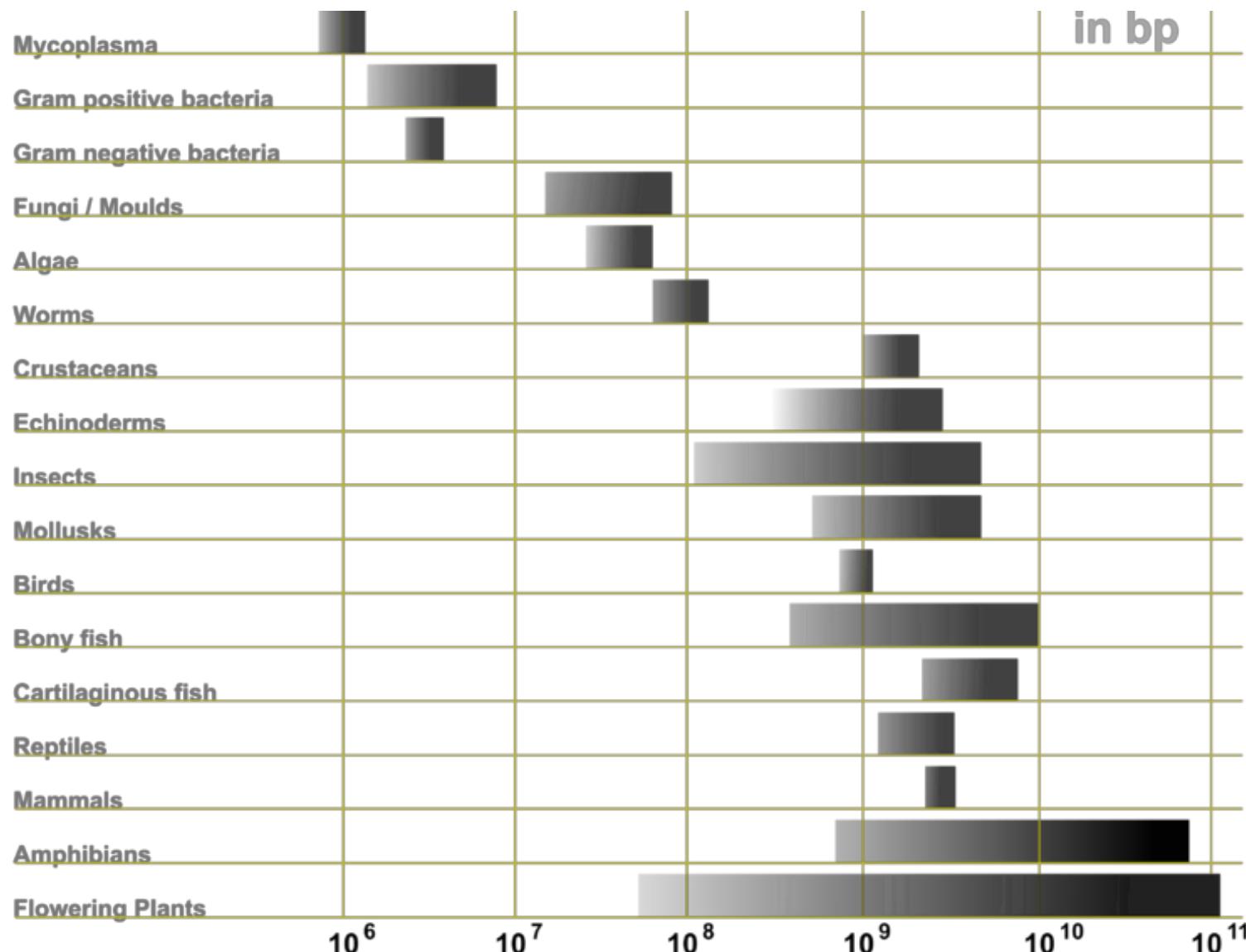


TABLE 1–2

A Few of the Many Organisms Whose Genomes Have Been Completely Sequenced

Organism	Genome size (millions of nucleotide pairs)	Number of genes	Biological interest
<i>Mycoplasma genitalium</i>	0.58	483	Smallest true organism
<i>Treponema pallidum</i>	1.1	1,039	Causes syphilis
<i>Borrelia burgdorferi</i>	1.44	1,738	Causes Lyme disease
<i>Helicobacter pylori</i>	1.7	1,589	Causes gastric ulcers
<i>Methanococcus jannaschii</i>	1.7	1,783	Archaeon; grows at 85 °C!
<i>Haemophilus influenzae</i>	1.8	1,738	Causes bacterial influenza
<i>Archaeoglobus fulgidus*</i>	2.2	—	High-temperature methanogen
<i>Synechocystis</i> sp.	3.6	4,003	Cyanobacterium
<i>Bacillus subtilis</i>	4.2	4,779	Common soil bacterium
<i>Escherichia coli</i>	4.6	4,377	Some strains cause toxic shock syndrome
<i>Saccharomyces cerevisiae</i>	12.5	5,770	Unicellular eukaryote
<i>Plasmodium falciparum</i>	23	5,268	Causes human malaria
<i>Caenorhabditis elegans</i>	100	19,400	Multicellular roundworm
<i>Anopheles gambiae</i>	278	13,700	Malaria vector
<i>Arabidopsis thaliana</i>	157	25,500	Model plant
<i>Oryza sativa</i>	390	37,500	Rice
<i>Drosophila melanogaster</i>	140	13,000	Laboratory fly ("fruit fly")
<i>Mus musculus domesticus</i>	2.4×10^3	25,000	Laboratory mouse
<i>Pan troglodytes</i>	2.4×10^3	25,000	Chimpanzee
<i>Homo sapiens</i>	2.9×10^3	25,000	Human

*The number of genes is not yet determined.

Table 1-2

Lehninger Principles of Biochemistry, Fifth Edition

© 2008 W.H. Freeman and Company

High-throughput sequencing technologies (aka Next-generation sequencing. NGS)



www.technologyreview.com
www.gatc-biotech.com/

Illumina sequencing

MiniSeq



MiSeq



NextSeq



HiSeq 4000



HiSeq X Ten



MAX OUTPUT
8 Gb
MAX READ NUMBER
25 million
MAX READ LENGTH
2x150 bp

MAX OUTPUT
15 Gb
MAX READ NUMBER
25 million
MAX READ LENGTH
2x300 bp

MAX OUTPUT
120 Gb
MAX READ NUMBER
400 million
MAX READ LENGTH
2x150 bp

MAX OUTPUT
1500 Gb
MAX READ NUMBER
5 billion
MAX READ LENGTH
2x150 bp

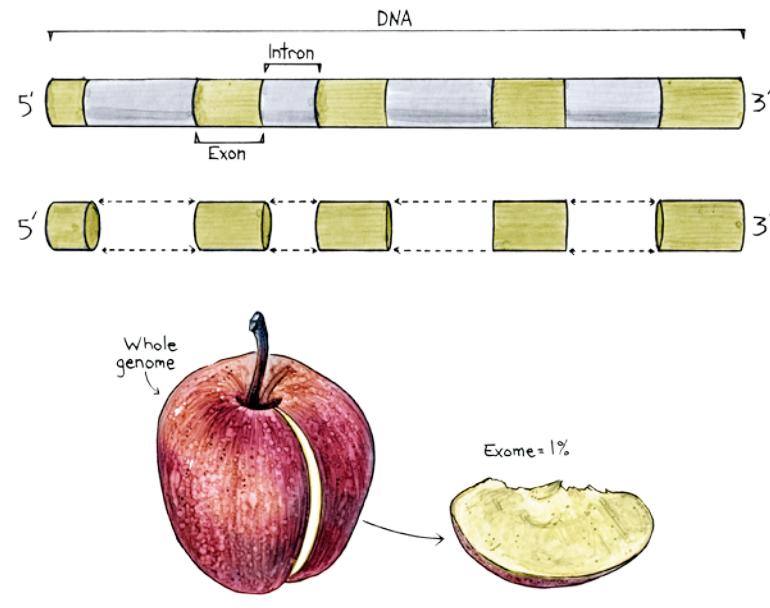
MAX OUTPUT
1800 Gb
MAX READ NUMBER
6 billion
MAX READ LENGTH
2x150 bp

1.5-2 TB of raw data per run (3 days)

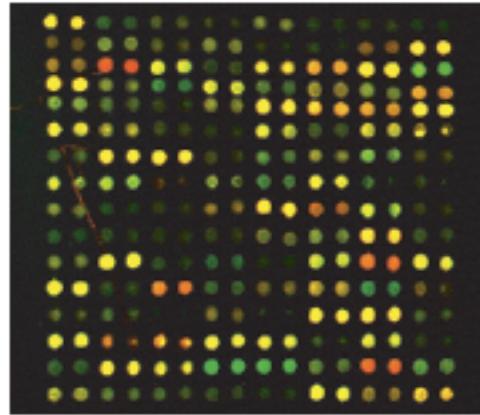
illumina®

Exome sequencing

- Key: Enriched for protein-coding regions of the genome
- Humans have about 180,000 exons
- That is, about 1% of the human genome
(30 million base pairs)

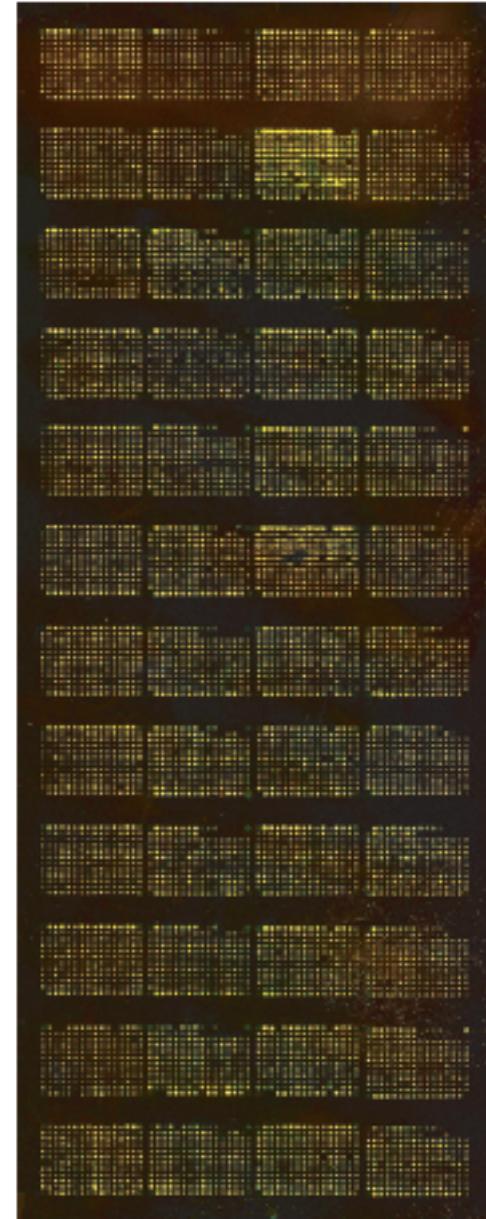


Microarray



Gene expression, genome variation, etc.

Typical data per chip: 10 - 100 MB



High-throughput proteomics

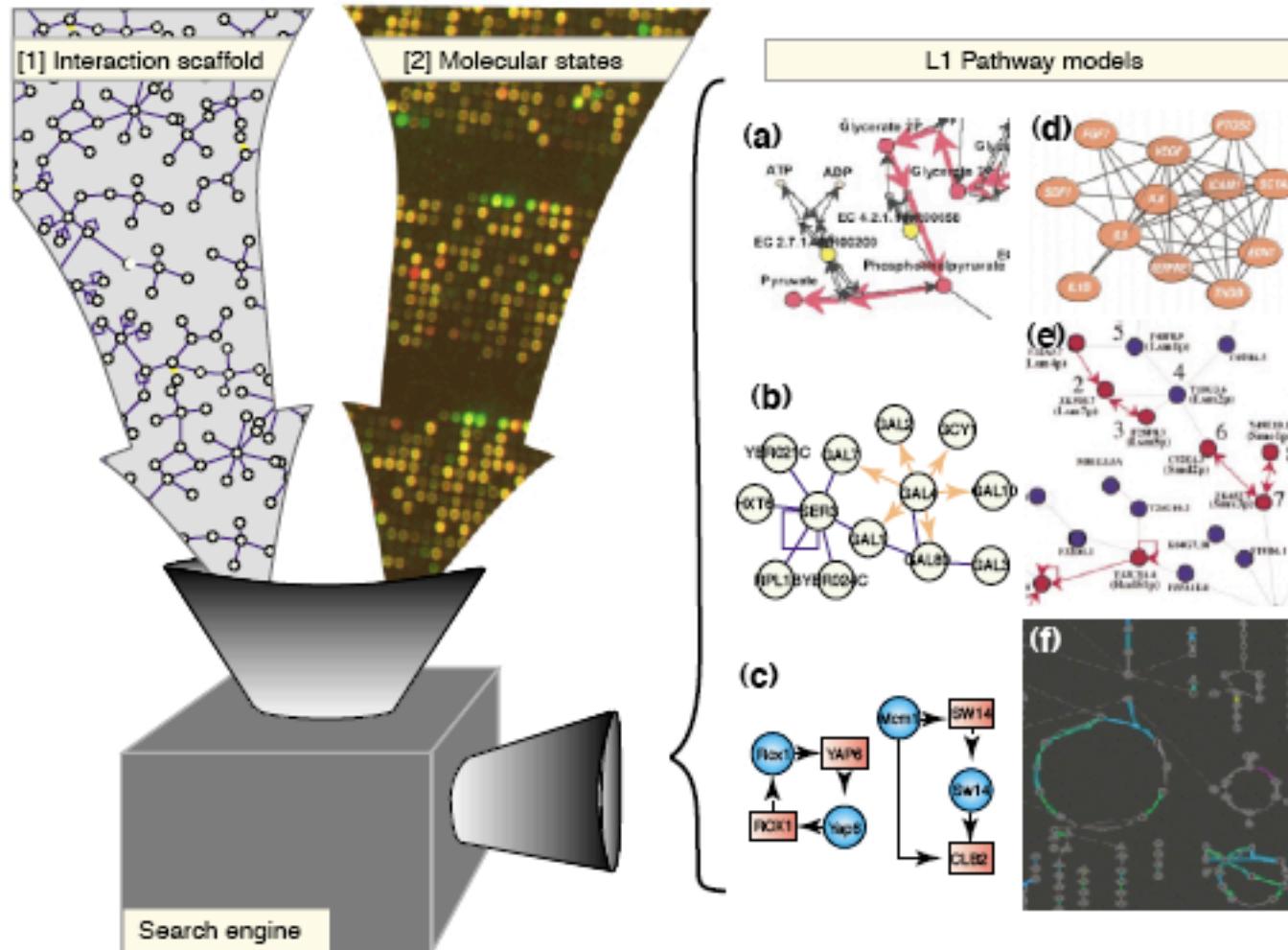


LC-MS/MS



LTQ Orbitrap

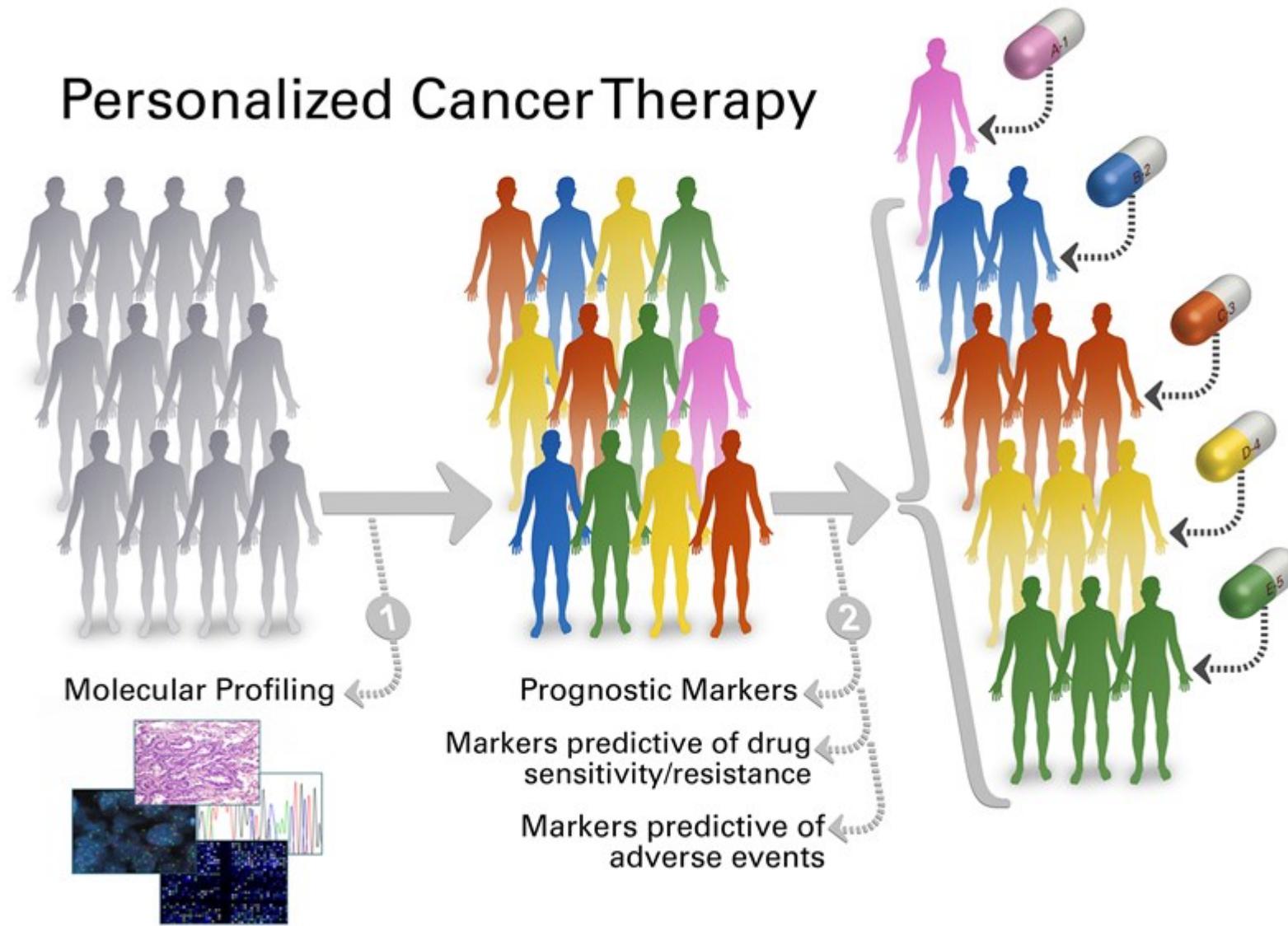
From large-scale data to meaningful findings



TRENDS in Biotechnology

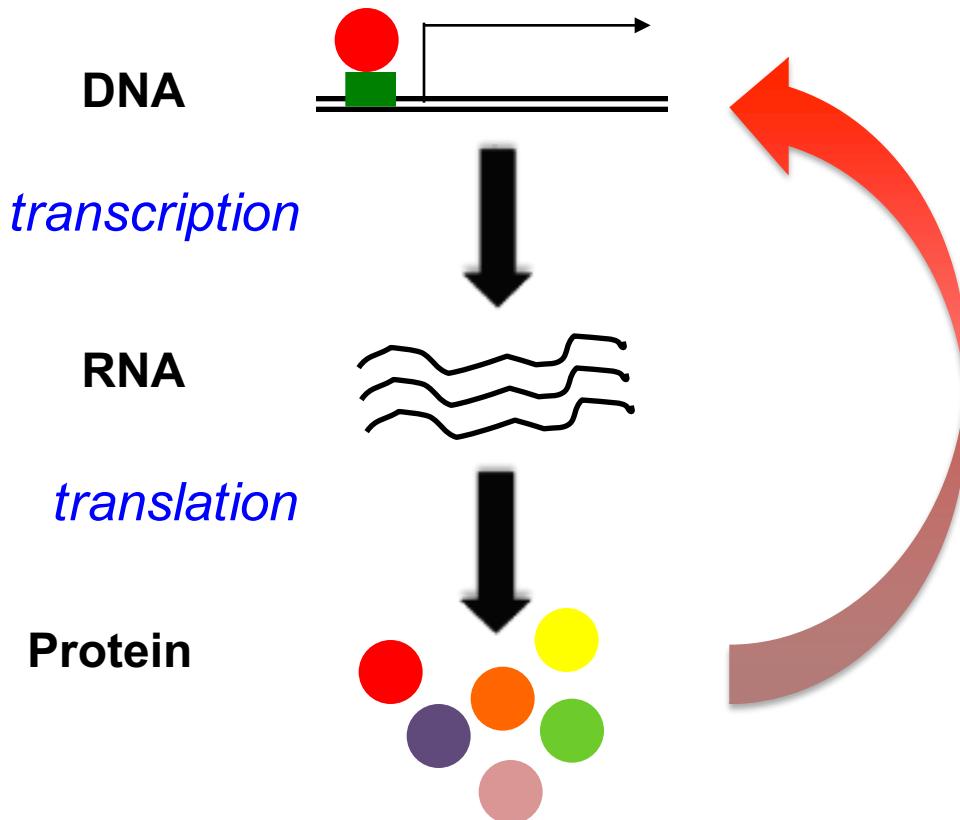
Ideker et al. Trends in Biotech (2003) Jun;21(6):255-62.

Personalized medicine?

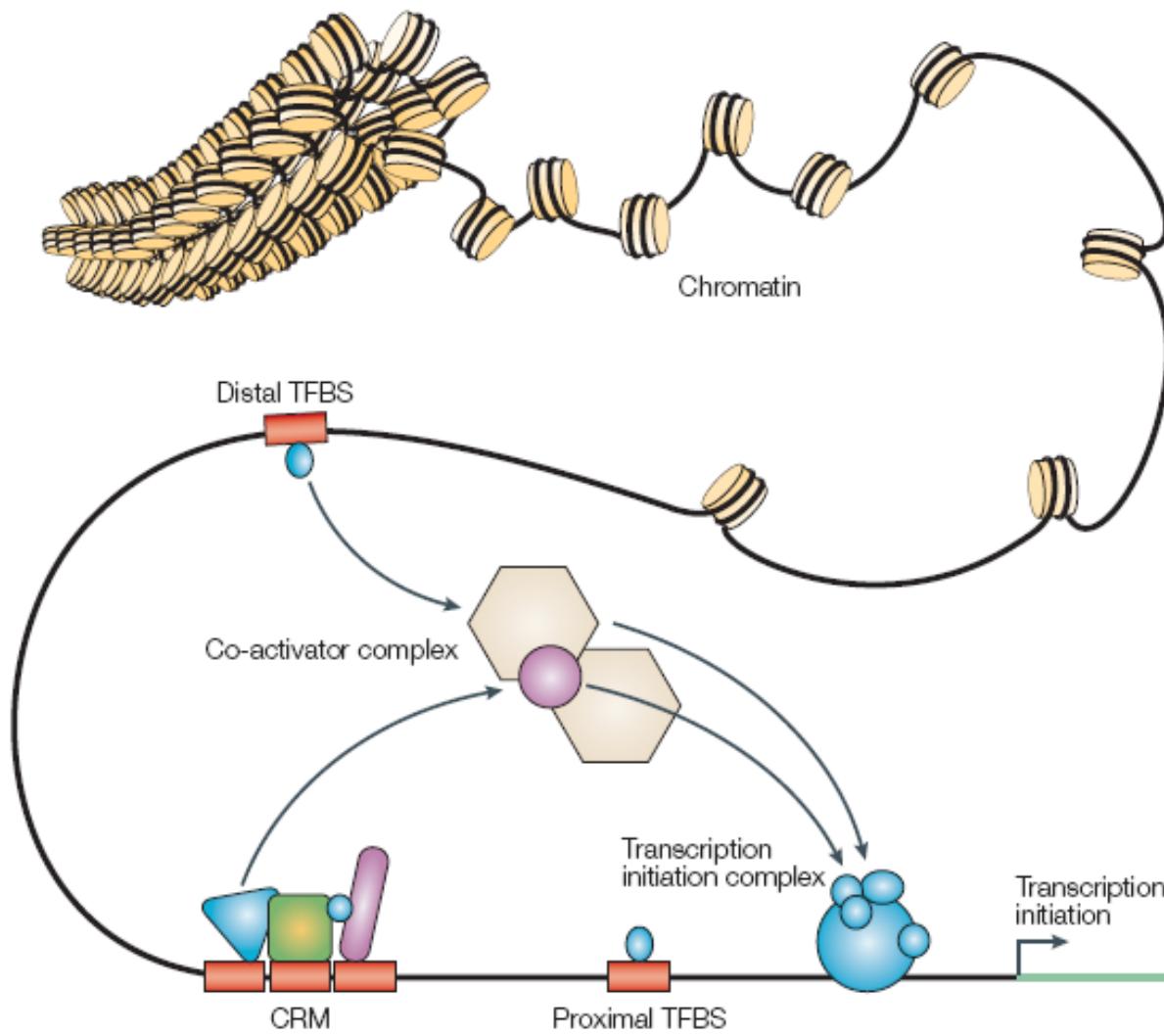


Pict from: gahtour.com/iran-looks-personalized-cancer-medicine-major-treatment/

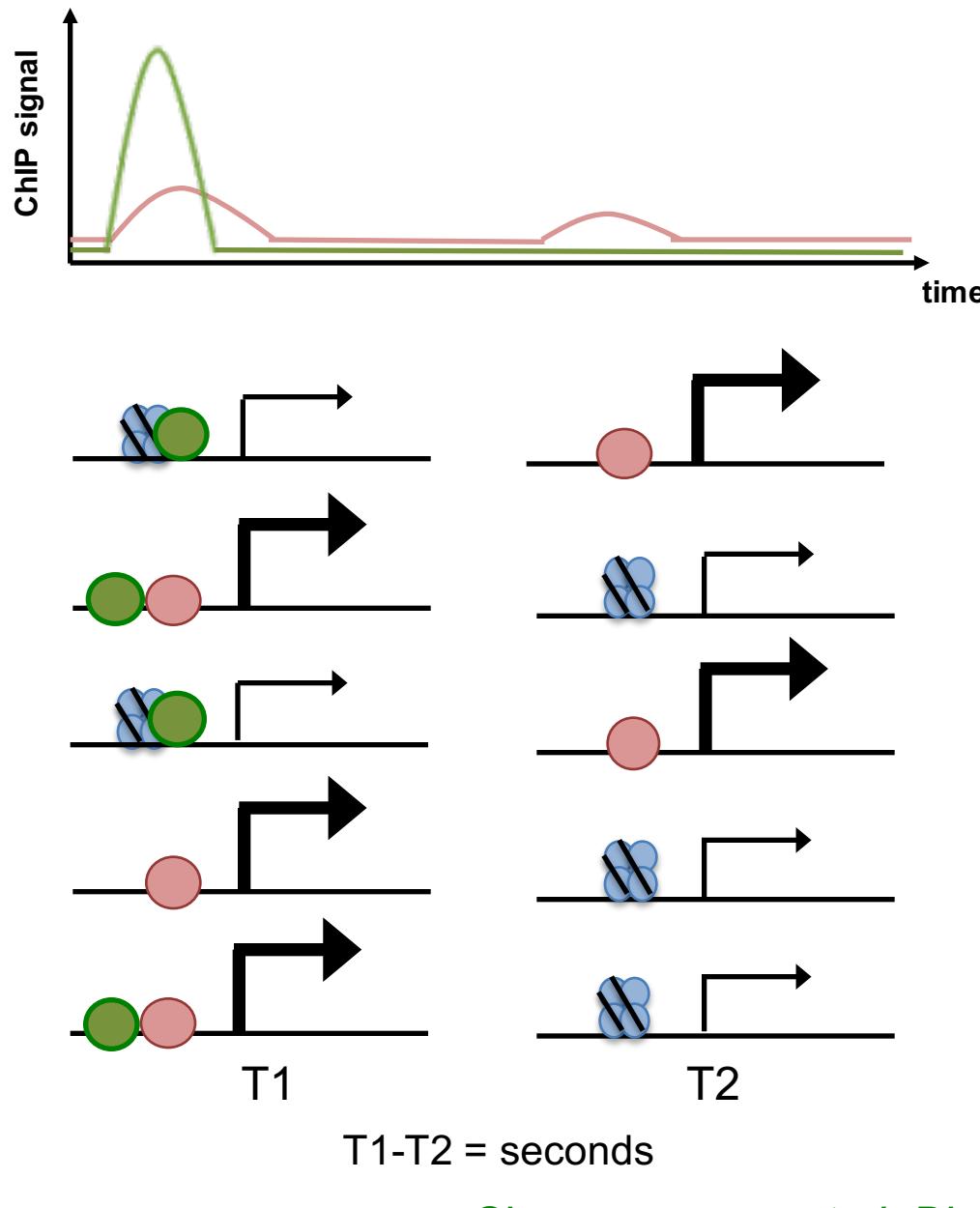
Central dogma of molecular biology



Transcription Factors Play a Central Role in Differential Expression of Genes



Transcriptional Regulation is highly dynamics



Transcriptional Regulation and Its Misregulation in Disease

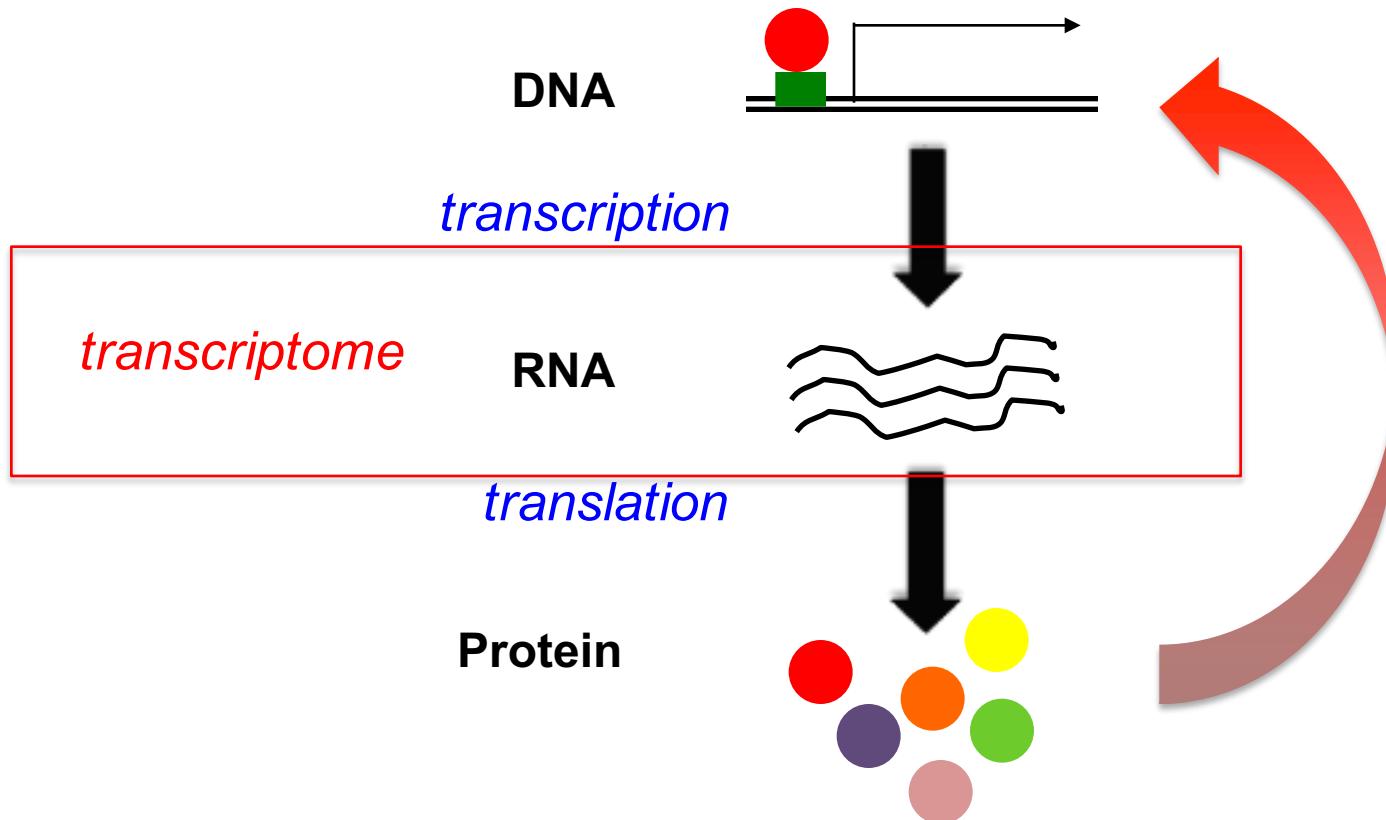
Tong Ihn Lee¹ and Richard A. Young^{1,2,*}

- **Various types of Cancers:**
 - NF-kappaB
 - STAT family proteins
 - P53
 - Steroid receptors
- **Autoimmunity and Inflammation**
- **Developmental Disorders**
- **Diabetes**
- **Cardiovascular Diseases**

Outlines

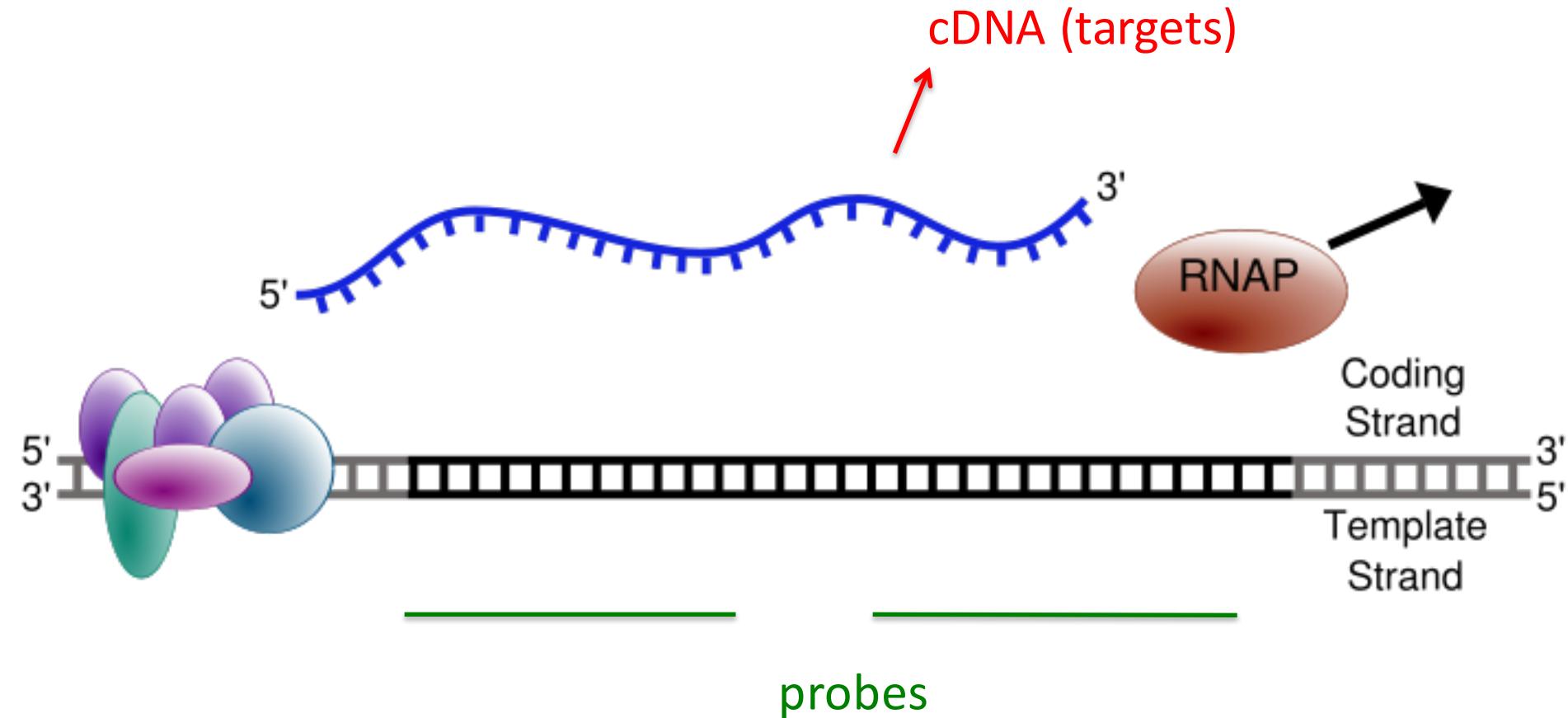
- Transcriptomics: genome-wide transcriptional read-outs
- Regulation of transcription (and transcriptomes)
- Putting things in perspectives

Central dogma of molecular biology

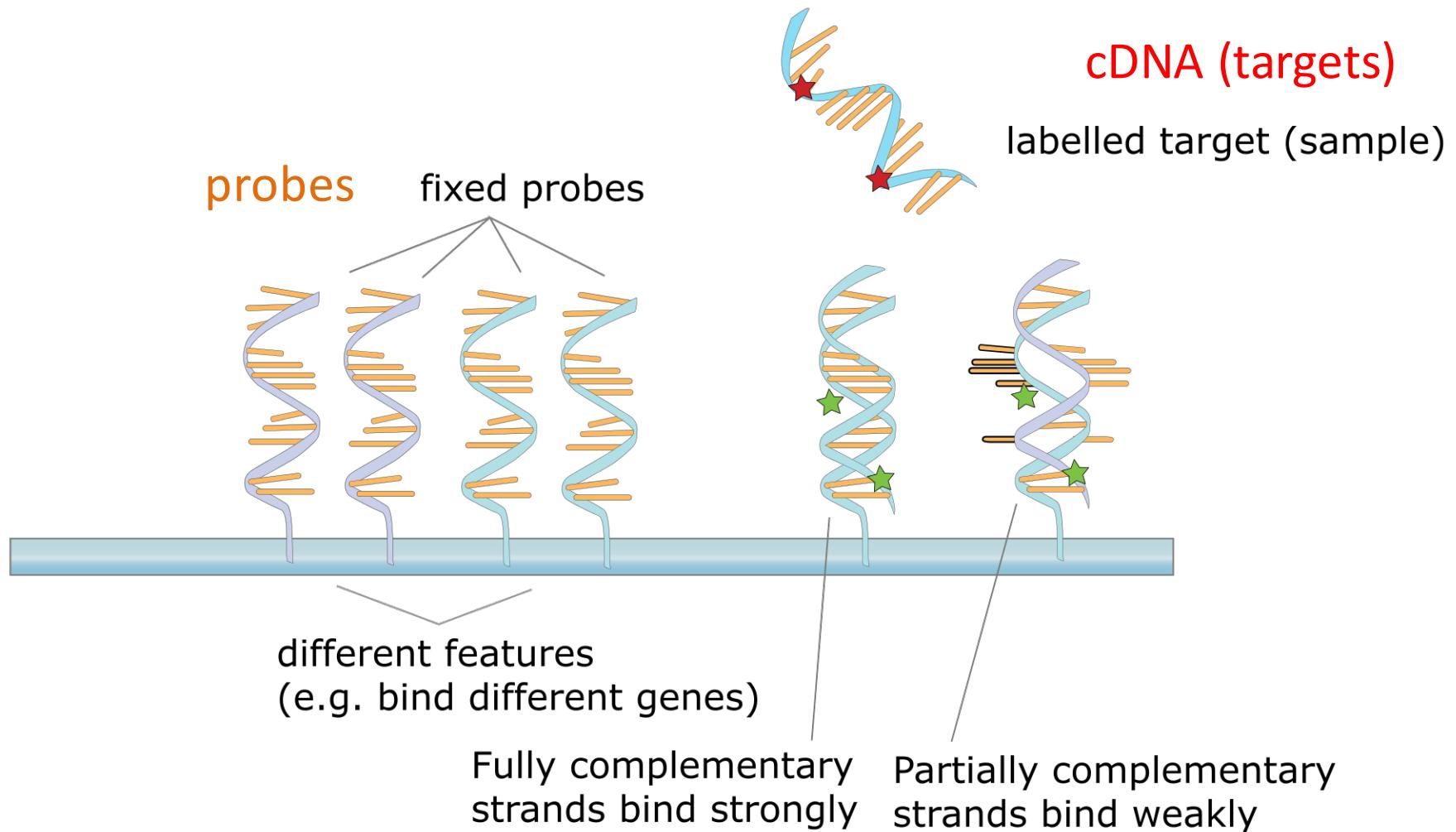


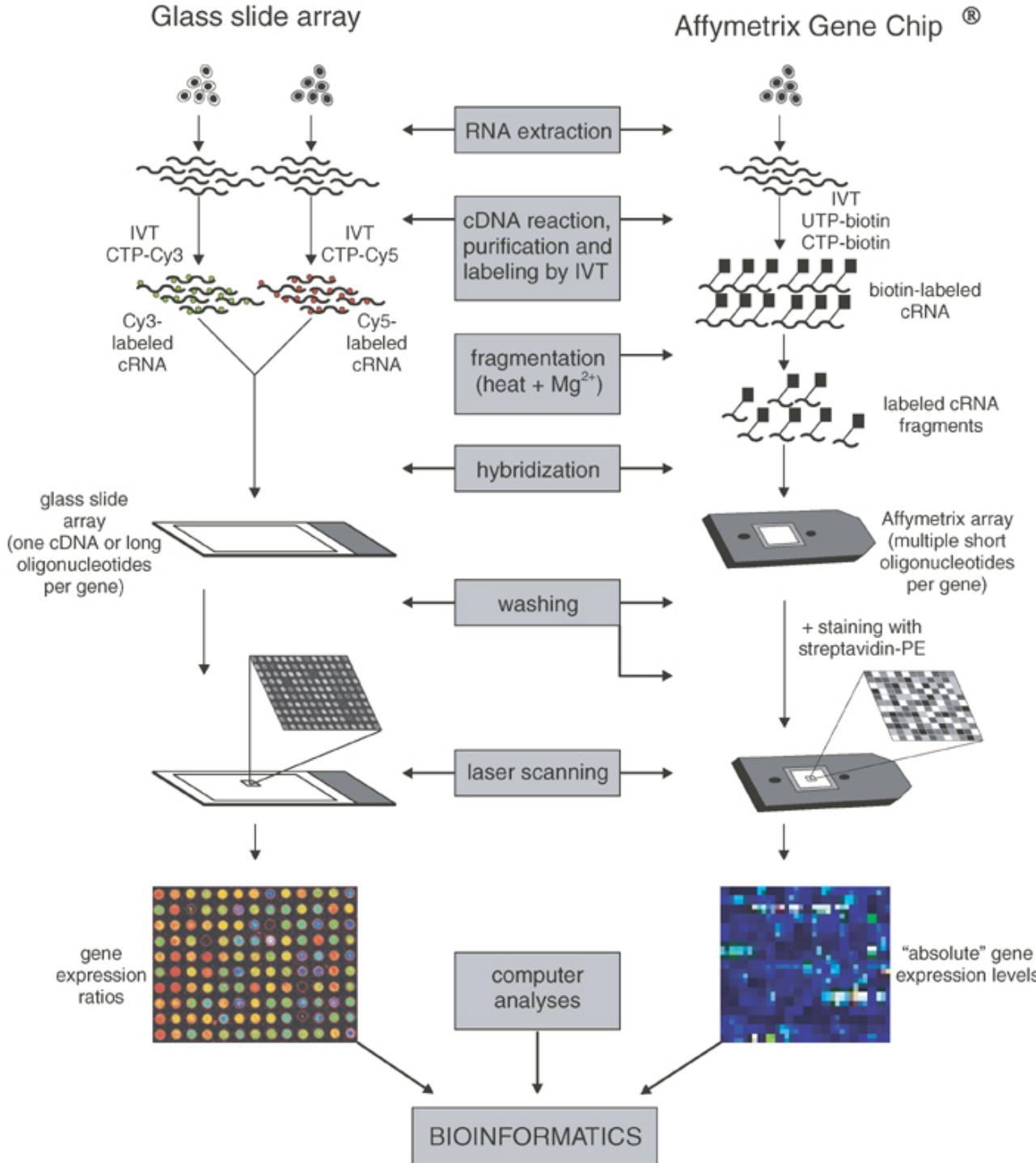
How do we determine transcriptional rates?

Expression Array / Gene Array

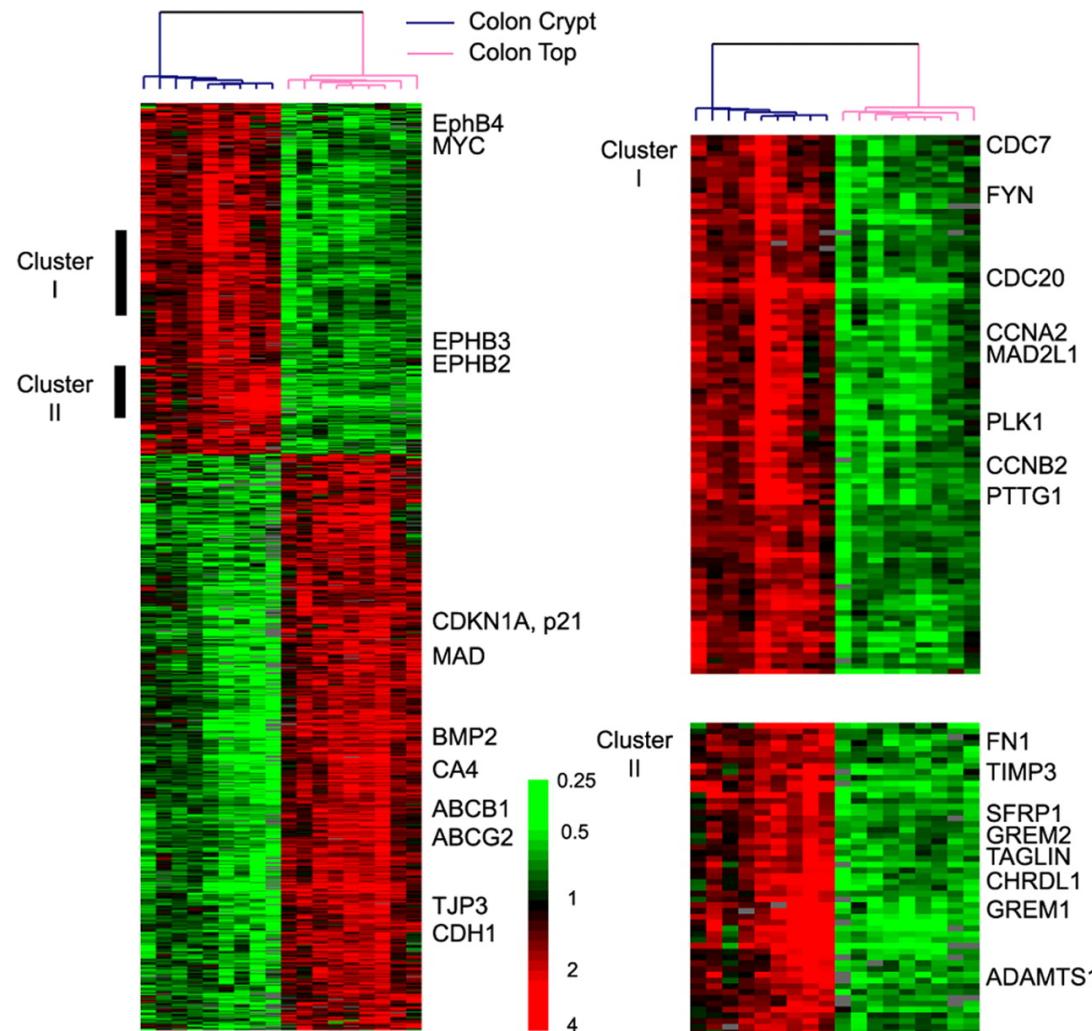


Microarray in Nutshell



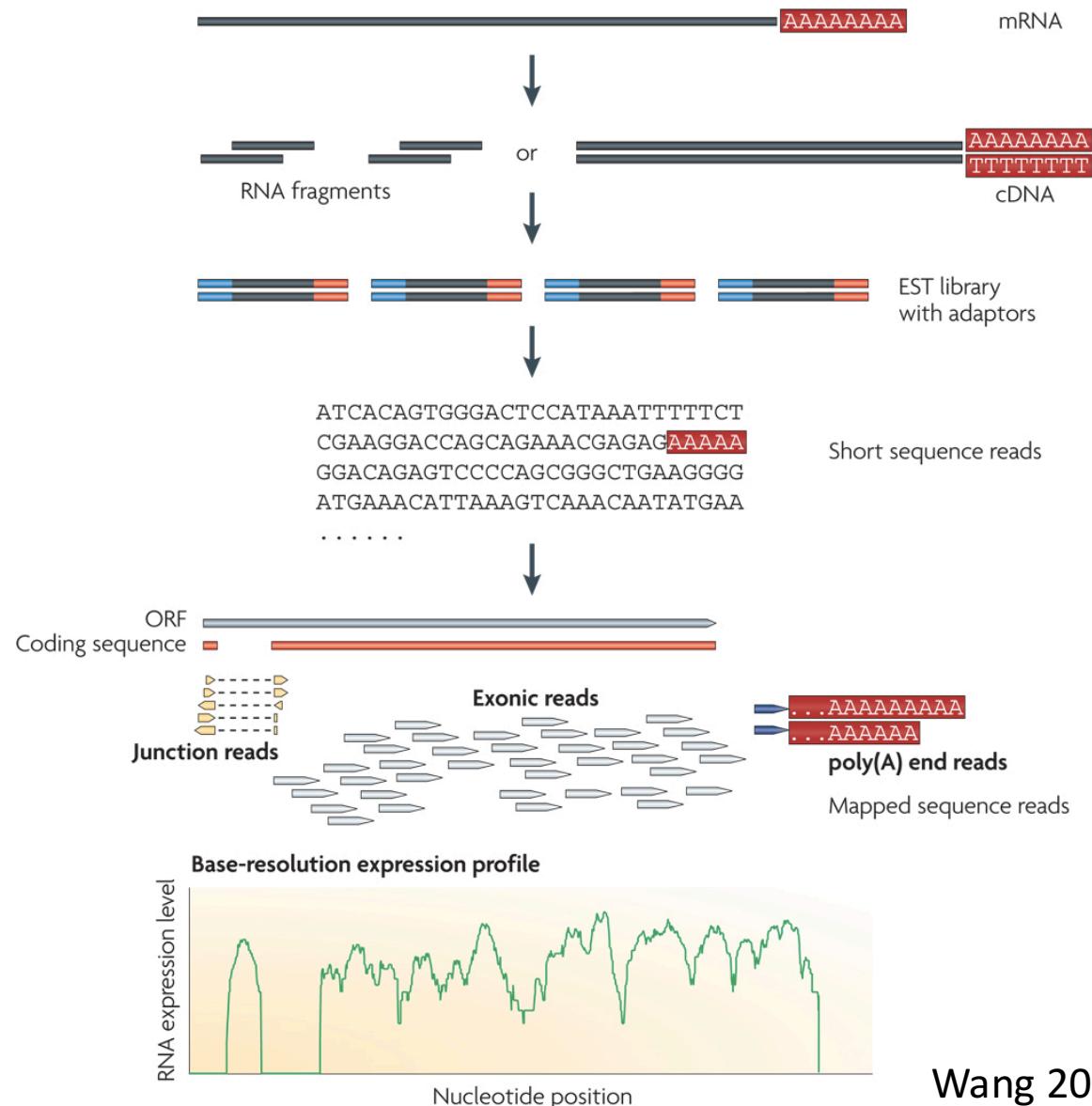


Discovery of new Biomarkers from Expression Array



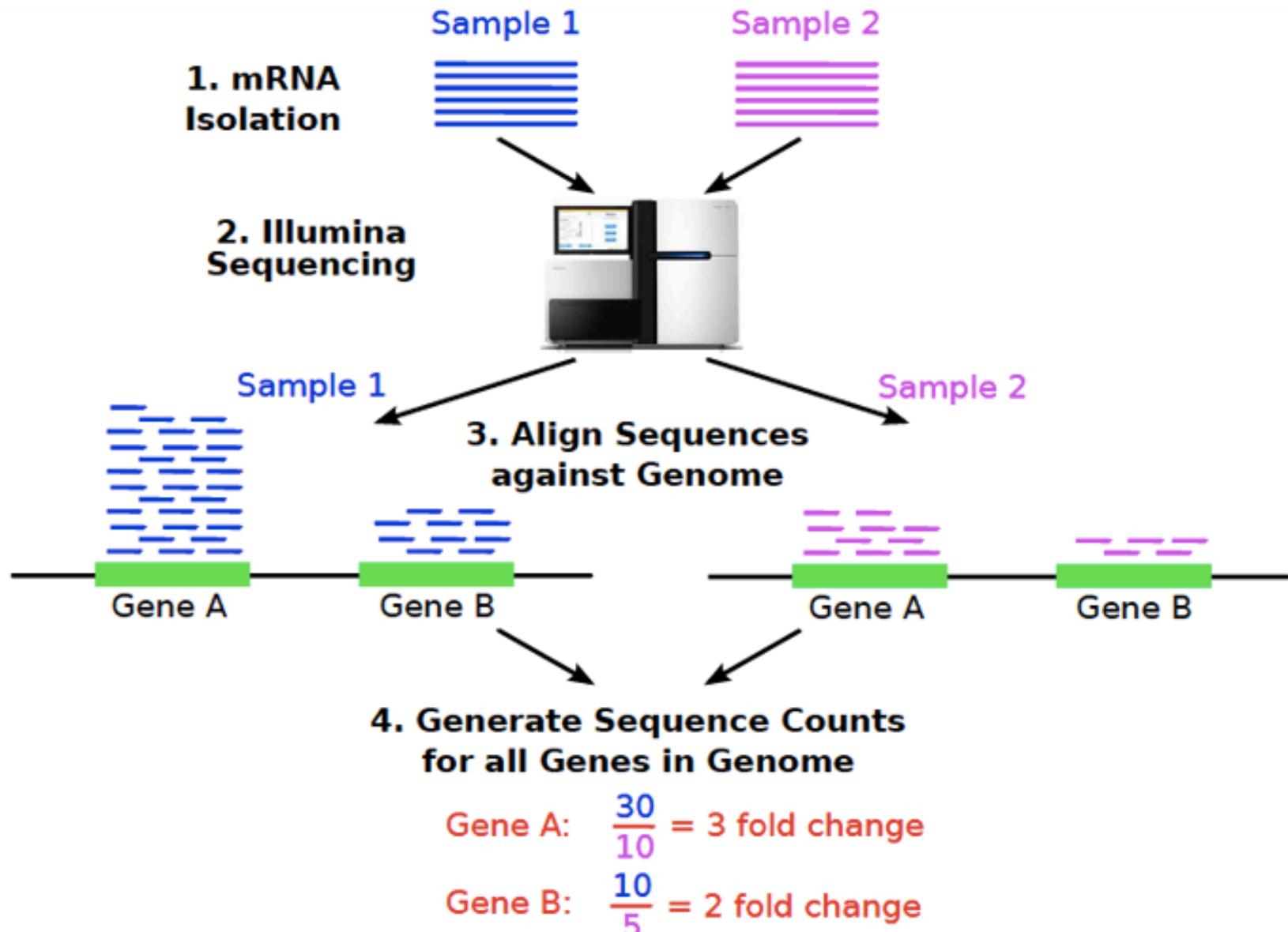
Cynthia Kosinski et al. PNAS 2007;104:15418-15423

RNA-seq: whole transcriptome shotgun sequencing

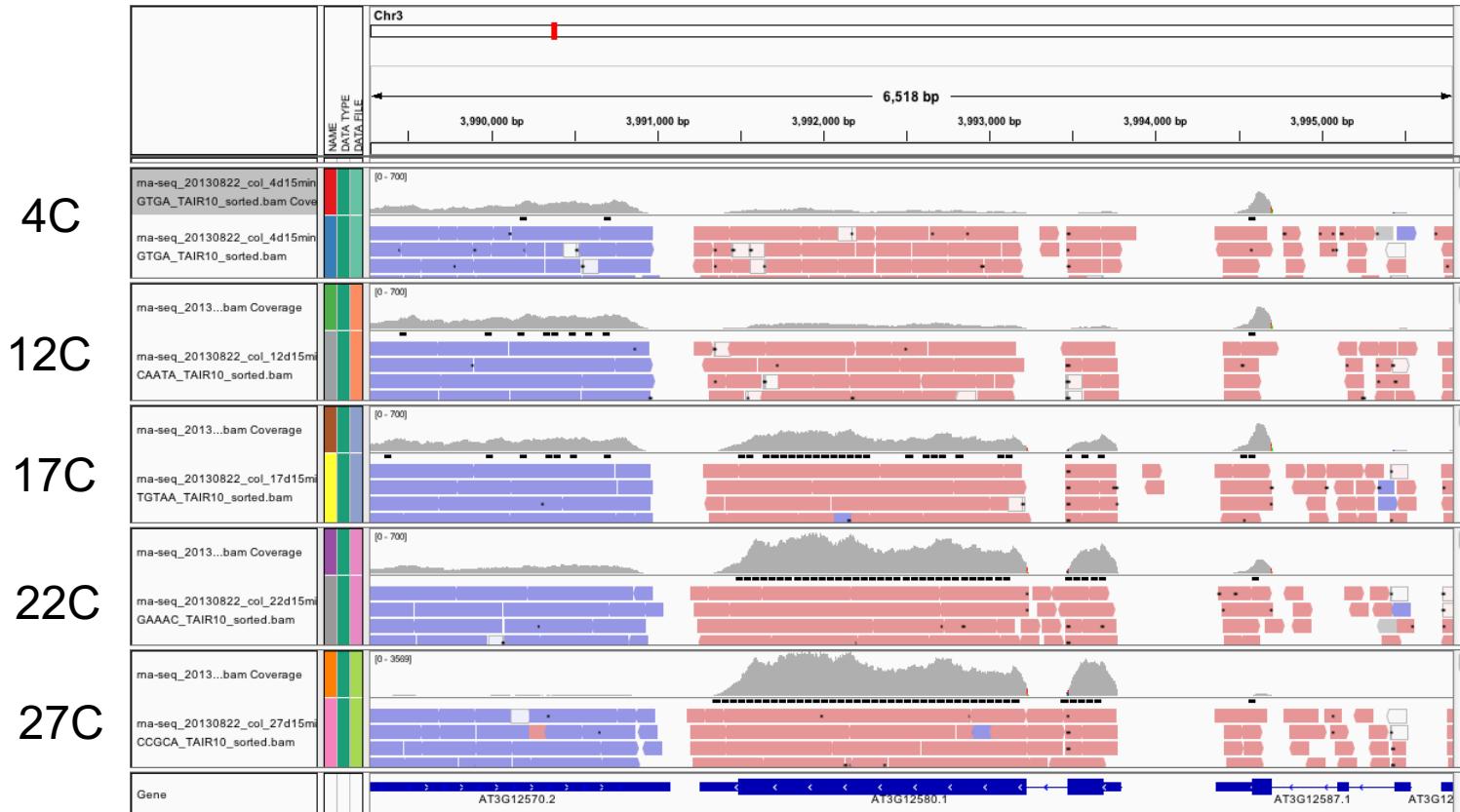


Wang 2009 *Nat Rev Genet.*

Differential Expressed mRNA



Examples of RNA-seq data: HSP70 as a temperature sensor



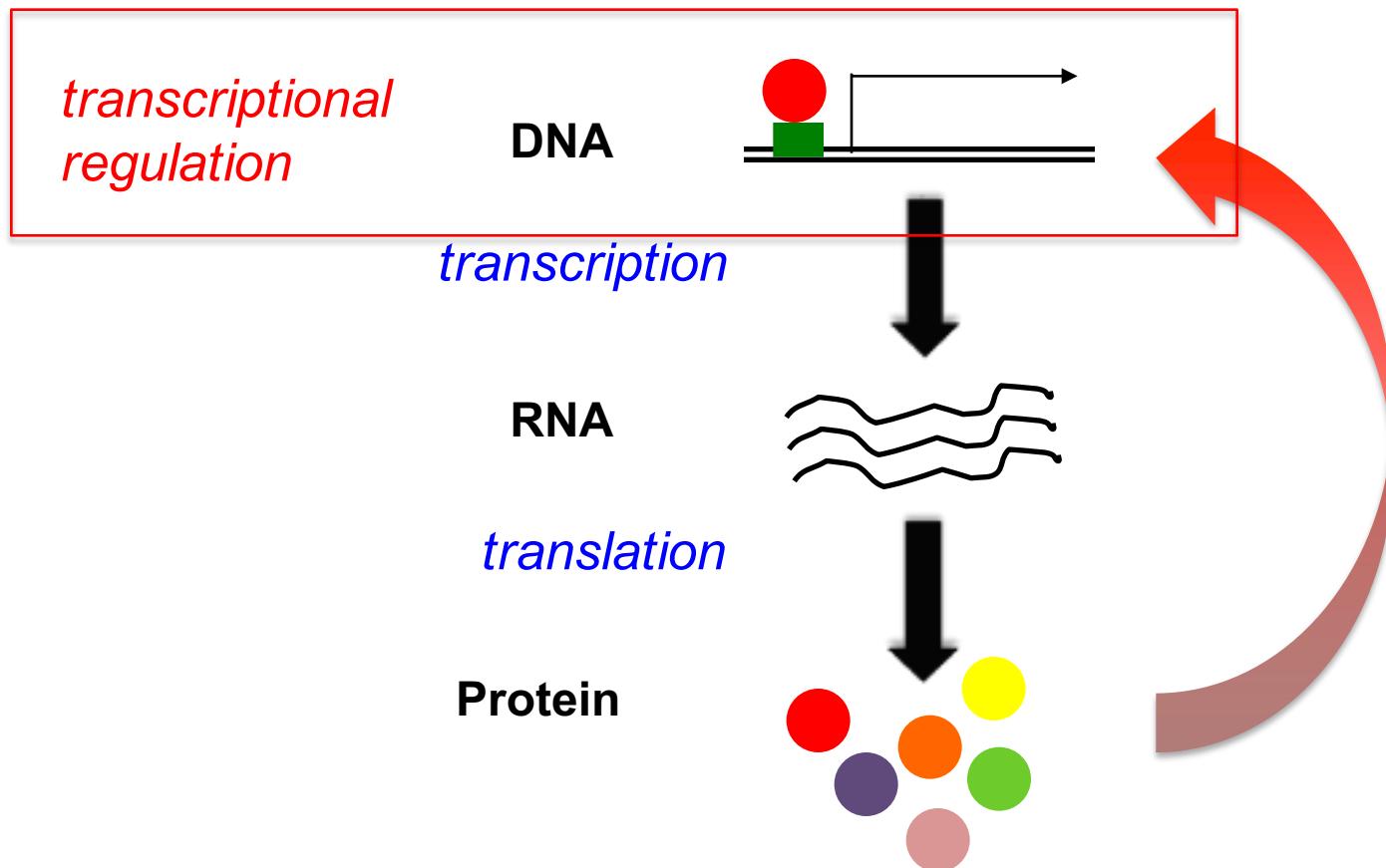
Transcriptional response after 15 mins

With Sandra Cortijo and Phil Wigge (Scale 0-600 for all, except 27C: 0-3500)

Outlines

- Transcriptomics: genome-wide transcriptional read-outs (**microarray, RNA-seq**)
- Regulation of transcription (and transcriptomes)
- Putting things in perspectives

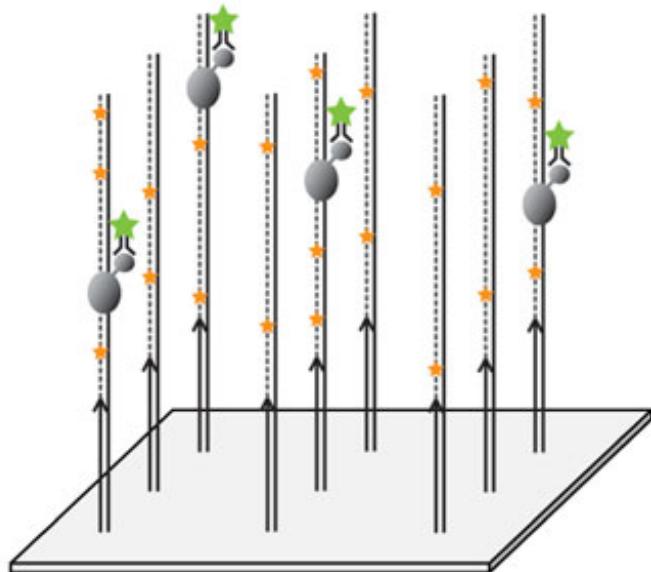
Central dogma of molecular biology



How to determine protein-DNA interaction?

in vitro

purified TF / naked DNA

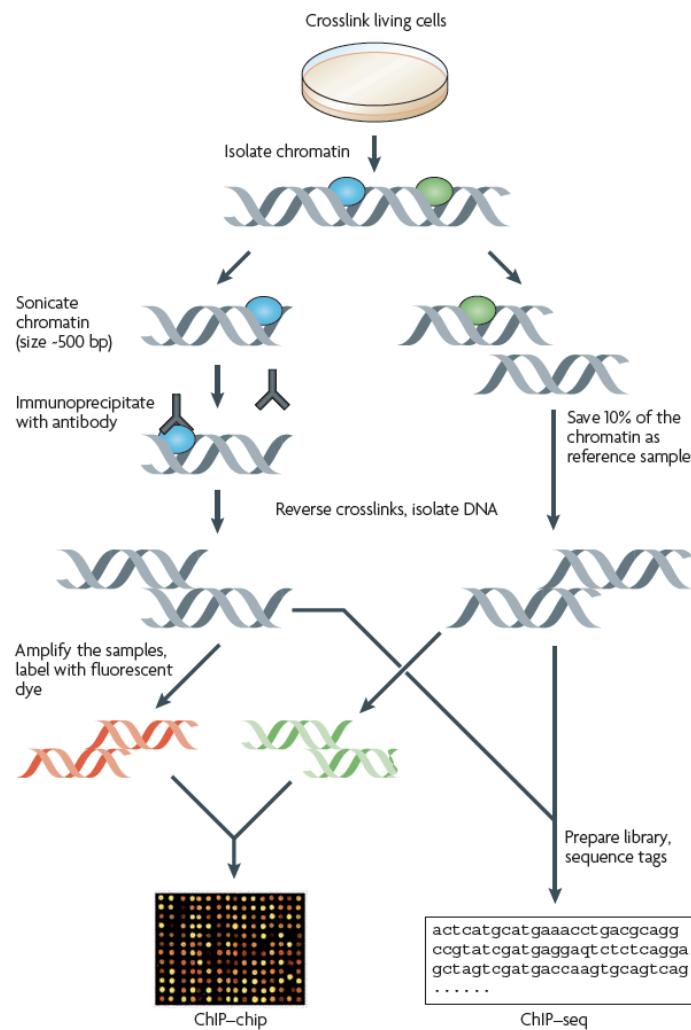


- ★ Cy3-labeled dUTP
- GST-tagged TF
- ★ Alexa488-labeled α -GST

e.g. Protein Binding Microarray, EMSA

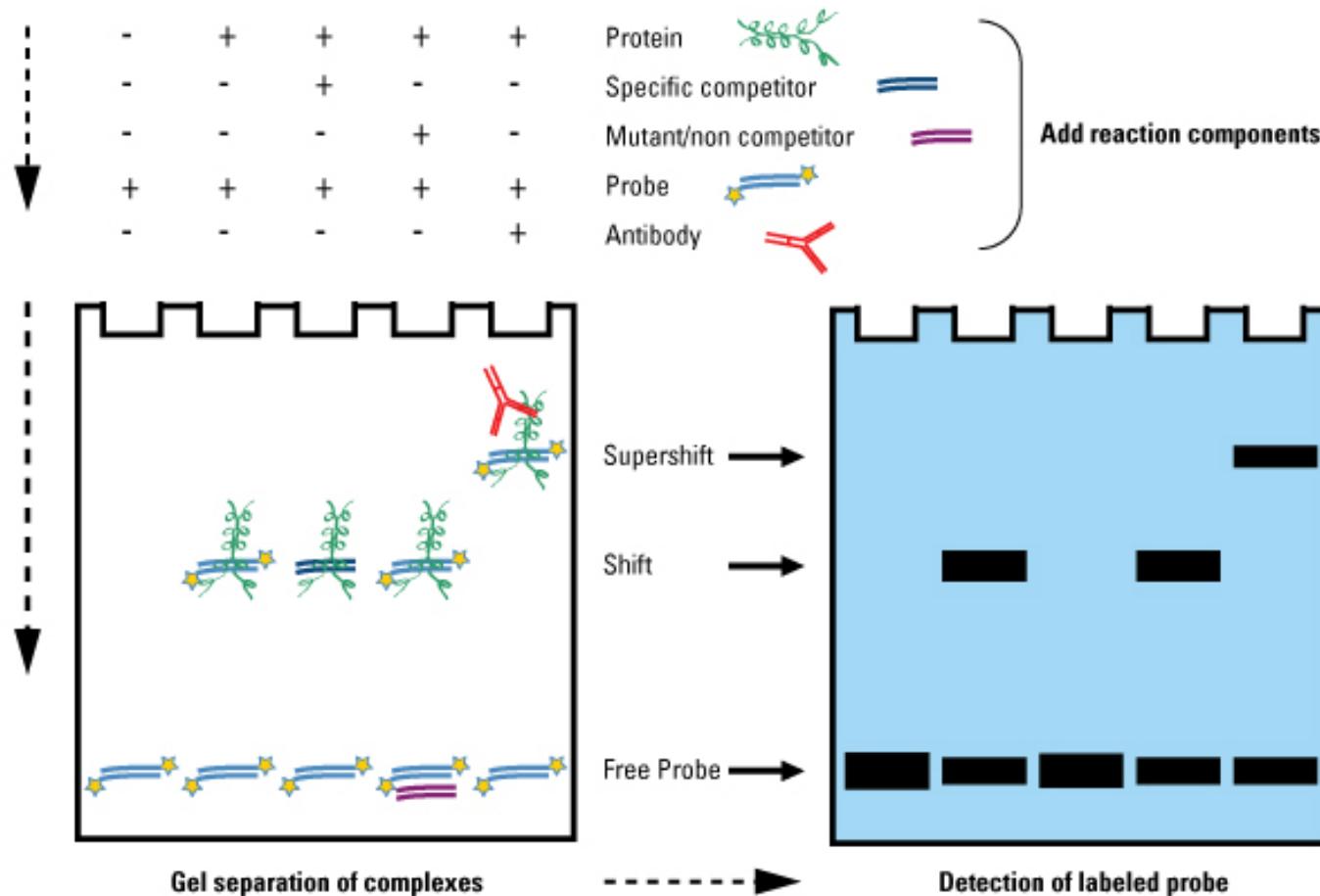
in vivo

all TFs / chromatin

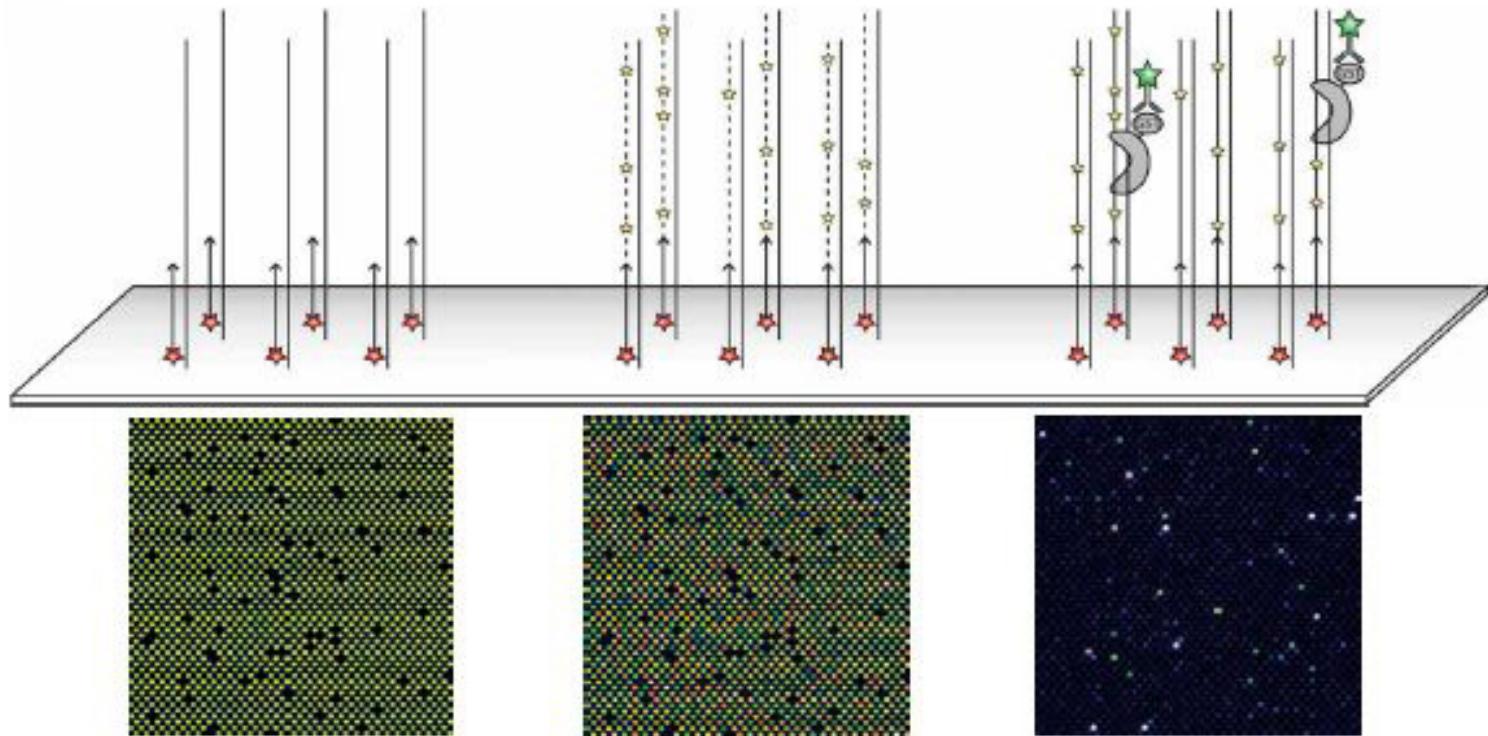


e.g. ChIP-chip, ChIP-seq

Electrophoretic mobility shift assay (EMSA)



Protein Binding Microarray (PBM)



http://the_brain.bwh.harvard.edu/

TF-DNA binding specificity (1)

Alignment of TF-bound sequences

Site 1	G	A	C	C	A	A	A	T	A	A	G	G	C	A
Site 2	G	A	C	C	A	A	A	T	A	A	G	G	C	A
Site 3	T	G	A	C	T	A	T	A	A	A	A	G	G	A
Site 4	T	G	A	C	T	A	T	A	A	A	A	G	G	A
Site 5	T	G	C	C	A	A	A	A	G	T	G	G	T	C
Site 6	C	A	A	C	T	A	T	C	T	T	G	G	G	C
Site 7	C	A	A	C	T	A	T	C	T	T	G	G	G	C
Site 8	C	T	C	C	T	T	A	C	A	T	G	G	G	C
	1	2	3	4	5	6	7	8	9	10	11	12	13	14
	Source binding sites													

: Position frequency matrix (PFM)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
A	0	4	4	0	3	7	4	3	5	4	2	0	0	4
C	3	0	4	8	0	0	0	3	0	0	0	0	2	4
G	2	3	0	0	0	0	0	0	1	0	6	8	5	0
T	3	1	0	0	5	1	4	2	2	4	0	0	1	0

Position weight matrix (PWM)

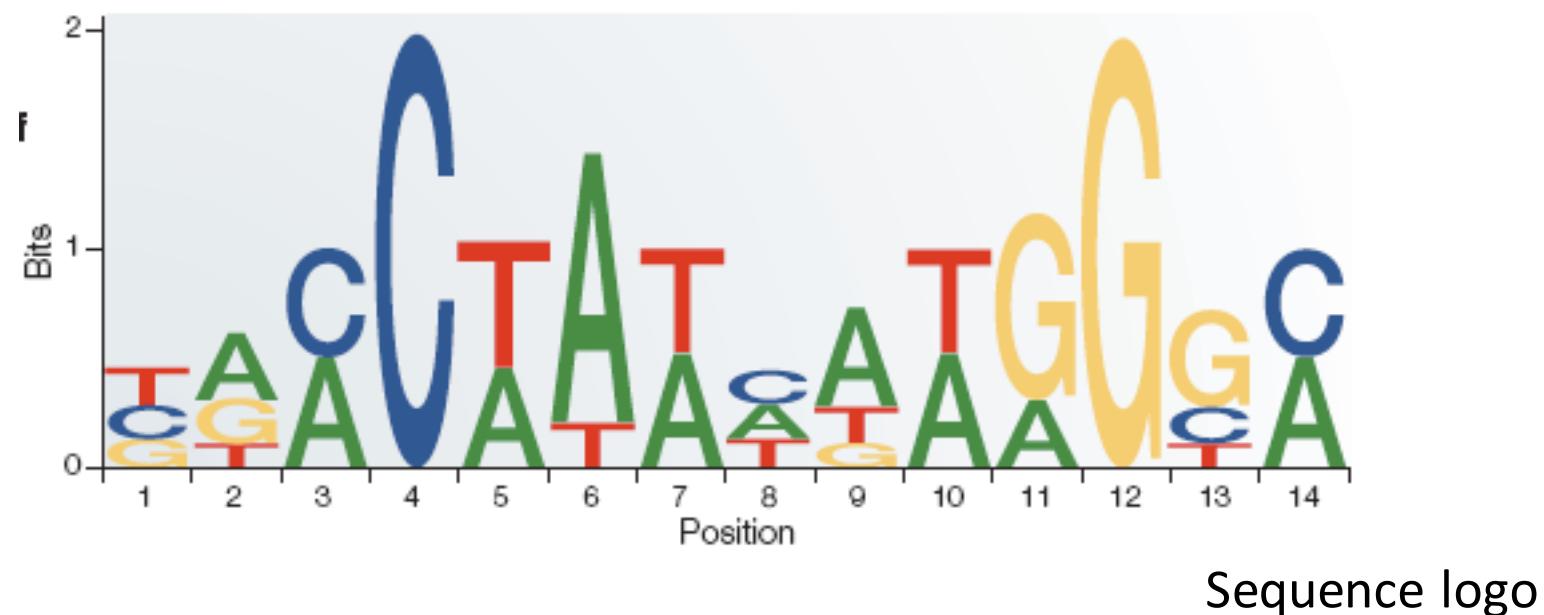
A	-1.93	0.79	0.79	-1.93	0.45	1.50	0.79	0.45	1.07	0.79	0.00	-1.93	-1.93	0.79
C	0.45	-1.93	0.79	1.68	-1.93	-1.93	-1.93	0.45	-1.93	-1.93	-1.93	-1.93	0.00	0.79
G	0.00	0.45	-1.93	-1.93	-1.93	-1.93	-1.93	-1.93	0.66	-1.93	1.30	1.68	1.07	-1.93
T	0.15	0.66	-1.93	-1.93	1.07	0.66	0.79	0.00	0.00	0.79	-1.93	-1.93	-0.66	-1.93

See Wasserman and Sanderlin, Nature review genetics, 2004 for more details

TF-DNA binding specificity (2)

Position frequency matrix (PFM)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
A	0	4	4	0	3	7	4	3	5	4	2	0	0	4
C	3	0	4	8	0	0	0	3	0	0	0	0	2	4
G	2	3	0	0	0	0	0	0	1	0	6	8	5	0
T	3	1	0	0	5	1	4	2	2	4	0	0	1	0



See Wasserman and Sanderlin, Nature review genetics, 2004 for more details

Databases of regulatory sequences

TRANSFAC: <http://www.gene-regulation.com/pub/databases/transfac/doc/toc.html>

RSAT: <http://rsat.ulb.ac.be/>

JASPAR: jaspar.cgb.ki.se/

TFSEARCH:

<http://www.cbrc.jp/research/db/TFSEARCH.html>

PROMO: http://alggen.lsi.upc.es/cgi-bin/promo_v3/promo/promoinit.cgi?dirDB=TF_8.3

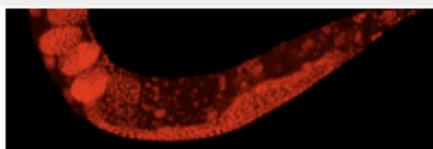


The high-quality transcription factor binding profile database

Browse the JASPAR CORE database directly:



JASPAR CORE Vertebrata



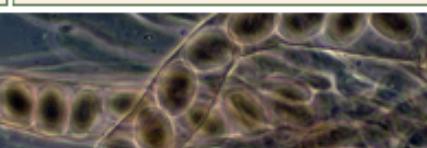
JASPAR CORE Nematoda



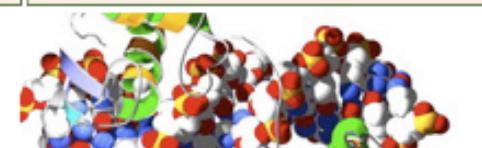
JASPAR CORE Insecta



JASPAR CORE Plantae



JASPAR CORE Fungi



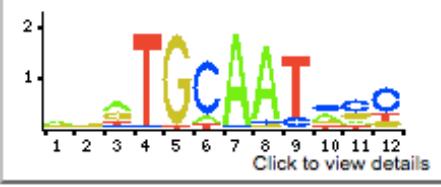
JASPAR CORE by Structural Class

[DOCUMENTATION](#)

[DOWNLOAD](#)

[CONTACT](#)

Finding potential regulator of gene X

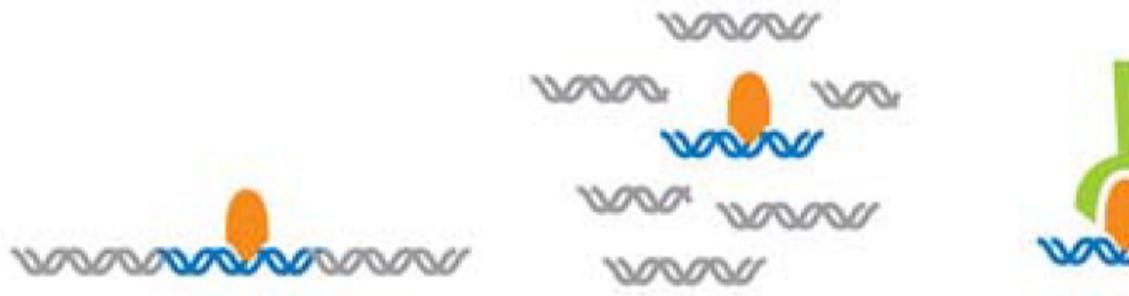
JASPAR matrix models:						
TOGGLE	ID	name	species	class	family	Sequence logo
<input checked="" type="checkbox"/>	MA0004.1	Arnt	Mus musculus	Basic helix-loop-helix factors (bHLH)	PAS domain factors	 Click to view details
<input checked="" type="checkbox"/>	MA0006.1	Ahr::Arnt	Mus musculus	Basic helix-loop-helix factors (bHLH)::Basic helix-loop-helix factors (bHLH)	PAS domain factors::PAS domain factors	 Click to view details
<input checked="" type="checkbox"/>	MA0019.1	Ddit3::Cebpa	Rattus norvegicus	Basic leucine zipper factors (bZIP)::Basic leucine zipper factors (bZIP)	C/EBP-related::C/EBP-related	 Click to view details
<input checked="" type="checkbox"/>	MA0025.1	NFIL3	Homo sapiens	Basic leucine zipper factors (bZIP)	C/EBP-related	 Click to view details

Finding potential regulator of gene X

11374 putative sites were predicted with these settings (80%) in sequence named 13

Model ID	Model name	Score	Relative score	Start	End	Strand	predicted site sequence
MA0476.1	FOS	1.335	0.806743455751239	1	11	-1	AATTAGACATC
MA0491.1	JUND	1.005	0.823280060697749	1	11	1	GATGTCTAATT
MA0775.1	MEIS3	6.075	0.868337860898837	1	8	-1	TAGACATC
MA0655.1	JDP2	5.484	0.8197004202266	2	10	1	ATGTCTAAT
MA0498.2	MEIS1	6.198	0.912644042924807	2	8	-1	TAGACAT
MA0899.1	HOXA10	6.405	0.836971511833828	4	14	-1	GGCAATTAGAC
MA0158.1	HOXA5	5.363	0.841814961828179	4	11	1	GTCTAATT
MA0909.1	HOXD13	4.488	0.804219559485583	4	13	-1	GCAATTAGAC
MA0913.1	Hoxd9	5.940	0.836166081301116	4	13	-1	GCAATTAGAC
MA0634.1	ALX3	8.462	0.897808355737197	5	14	-1	GGCAATTAGA
MA0634.1	ALX3	5.348	0.815185380813353	5	14	1	TCTAATTGCC
MA0151.1	Arid3a	5.399	0.826577864606885	5	10	-1	ATTAGA
MA0877.1	Barhl1	9.577	0.971239014736542	5	14	1	TCTAATTGCC
MA0635.1	BARHL2	6.366	0.889002993550392	5	14	1	TCTAATTGCC
MA0878.1	CDX1	5.650	0.84325666837775	5	13	-1	GCAATTAGA
MA0879.1	Dlx1	7.972	0.907268968995571	5	14	-1	GGCAATTAGA
MA0612.1	EMX1	5.726	0.824004796979857	5	14	-1	GGCAATTAGA
MA0612.1	EMX1	9.241	0.922767950639343	5	14	1	TCTAATTGCC

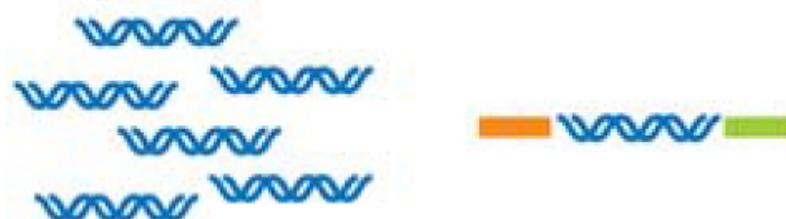
Chromatin immunoprecipitation (ChIP)



1. Cross-link bound proteins to DNA.

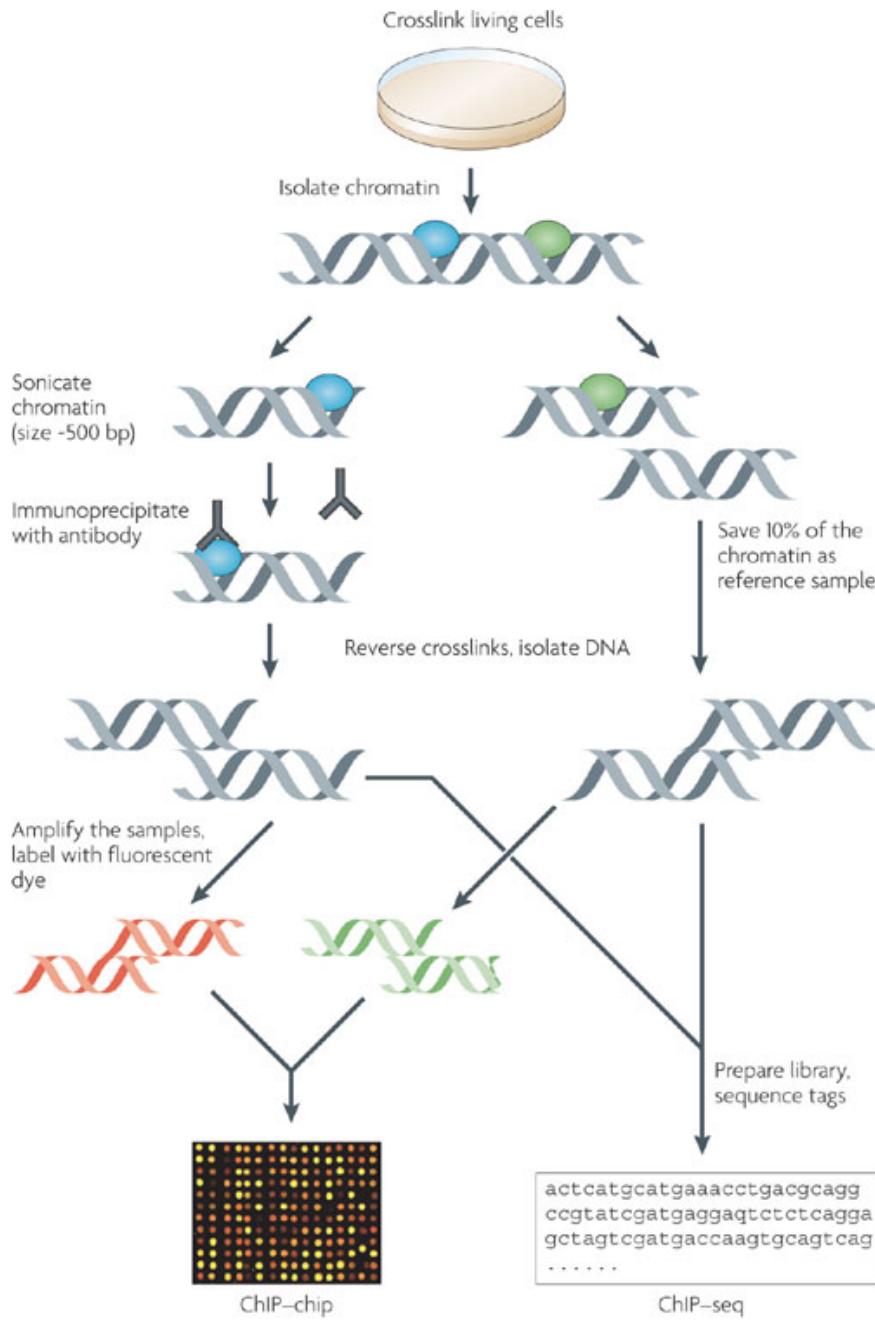
2. Isolate chromatin and shear DNA.

3. Precipitate chromatin with protein-specific antibody.

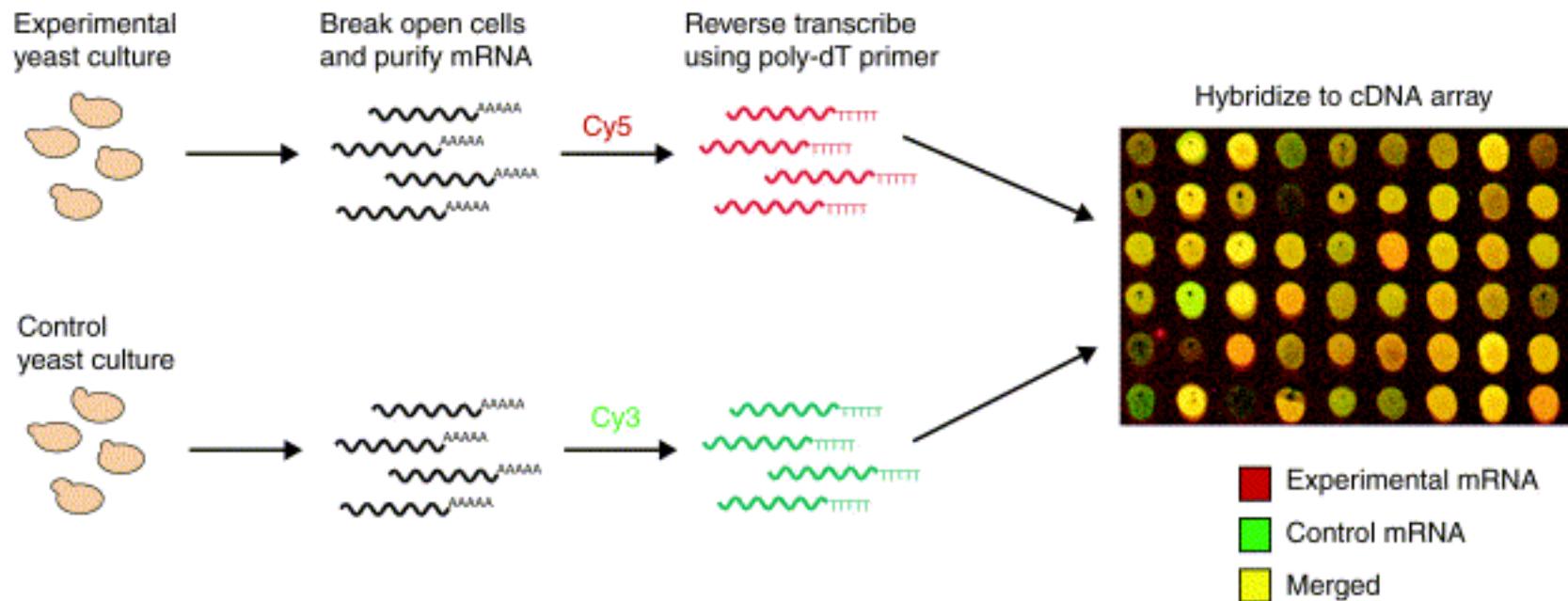


4. Reverse cross-link and digest protein.

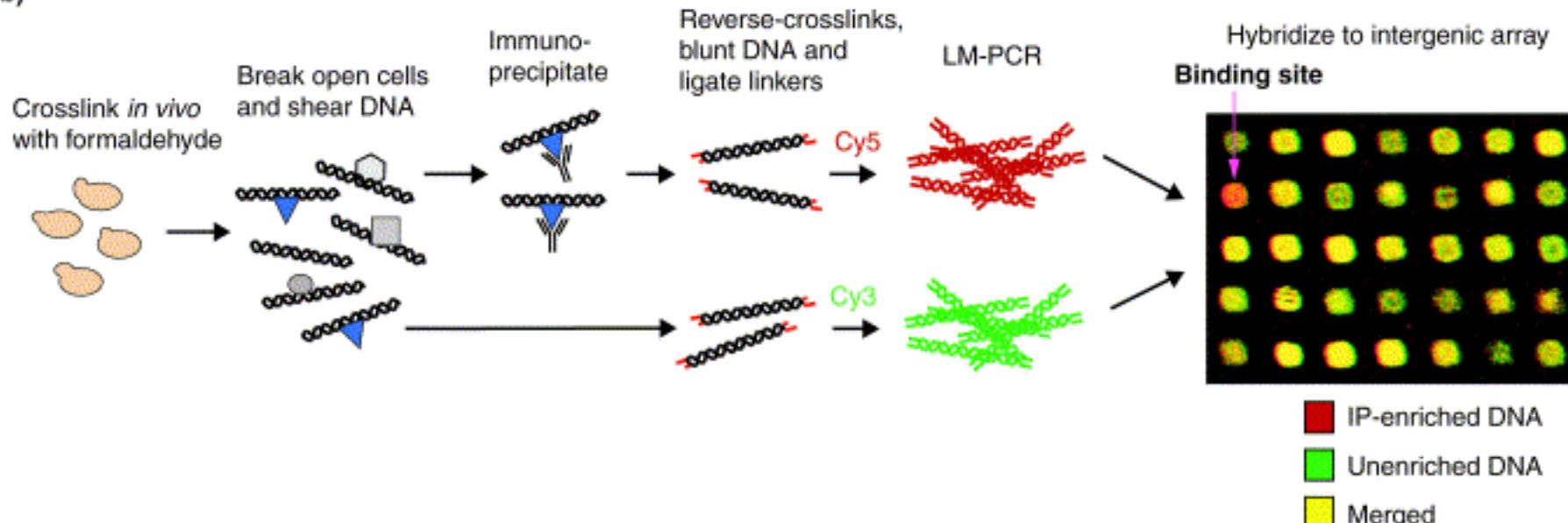
5. Ligate P1 and P2 adaptors to construct fragment library.



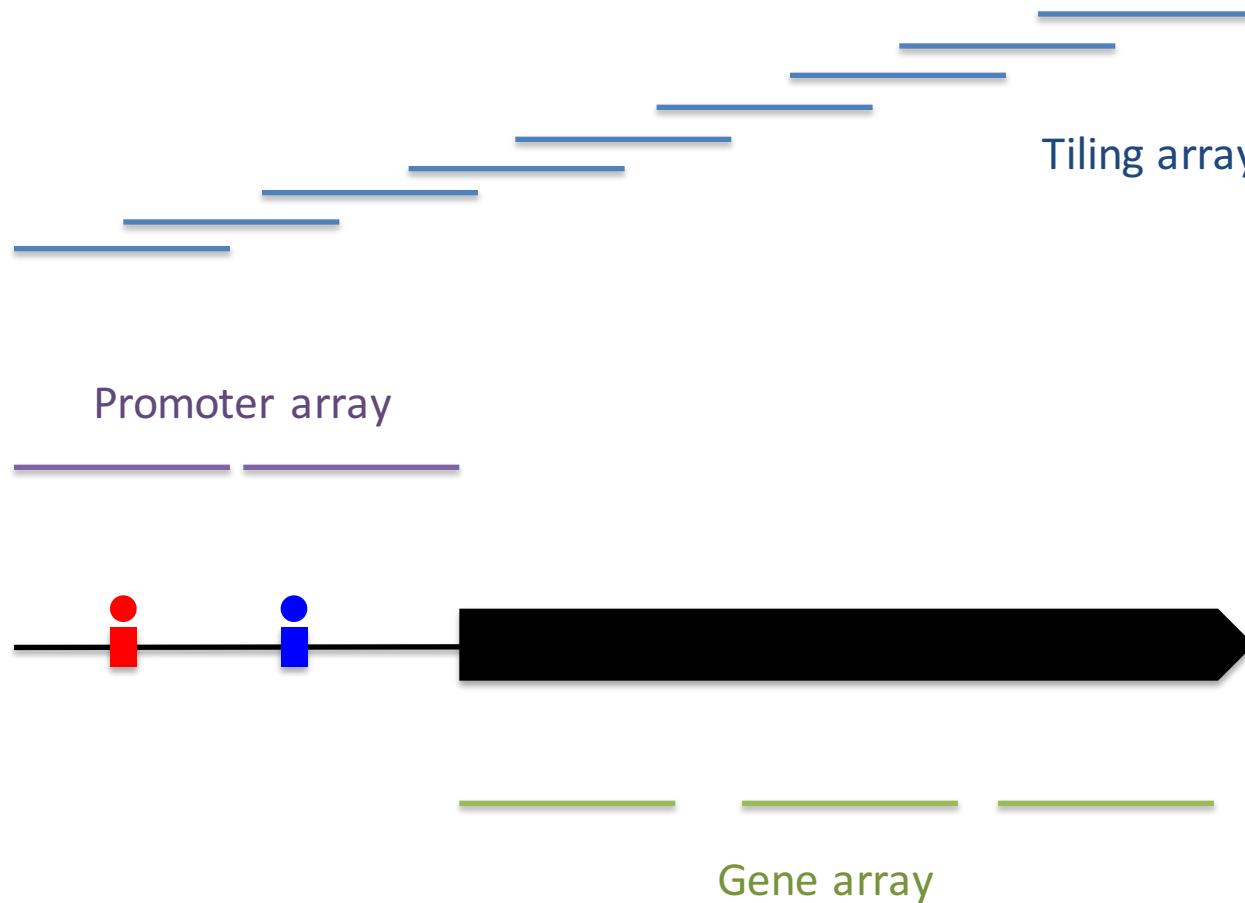
(a)



(b)

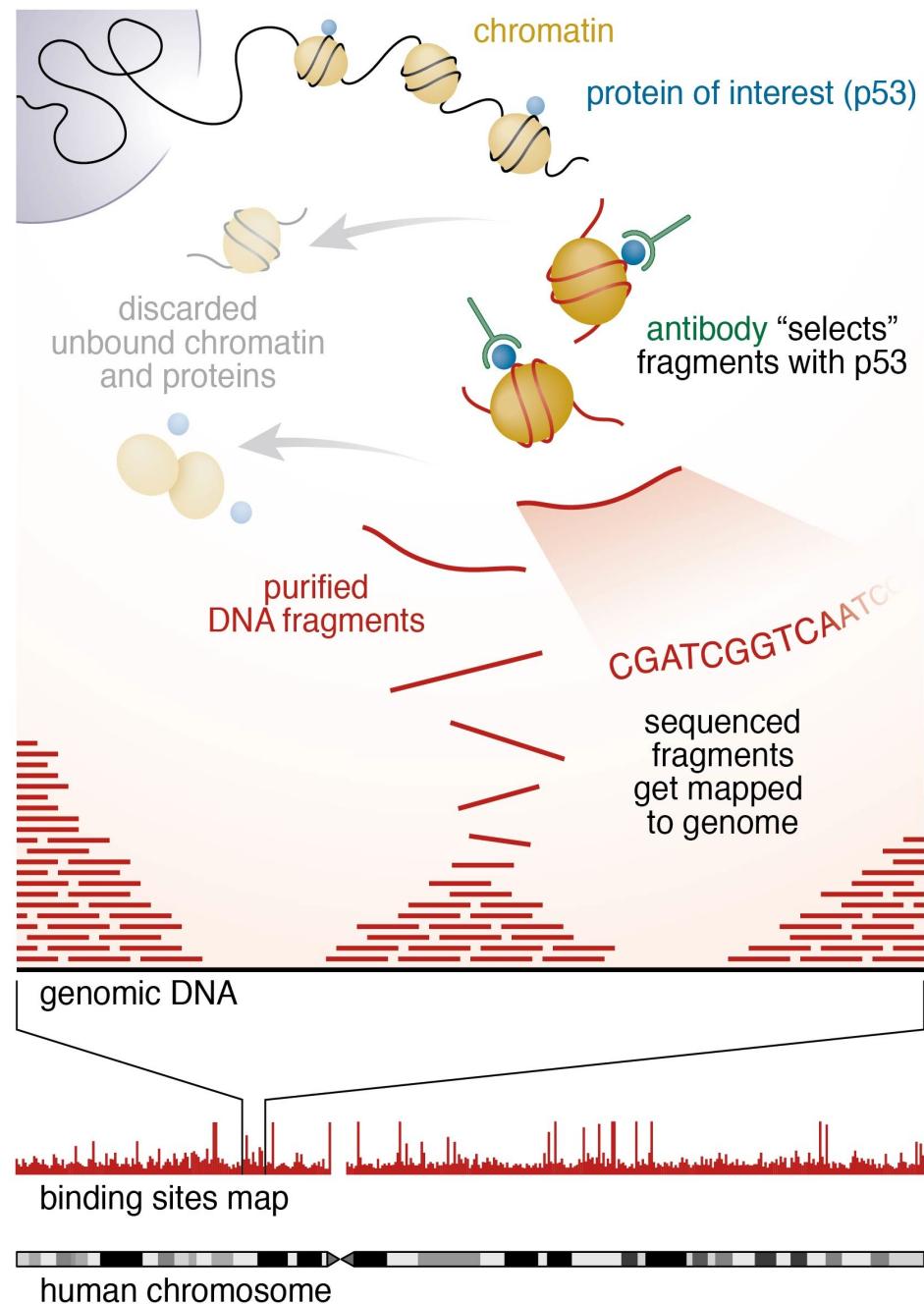


Arrays for ChIP-chip or Expression?



Chromatin

ImmunoPrecipitation – (next-generation) sequencing

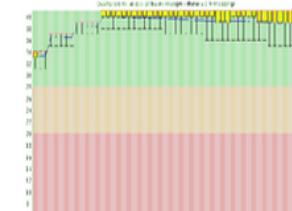


ChIP-seq data analysis

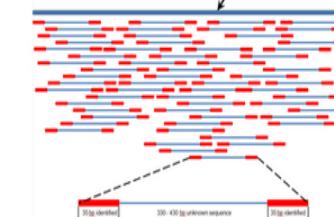
read data (FastQ)
rawdata



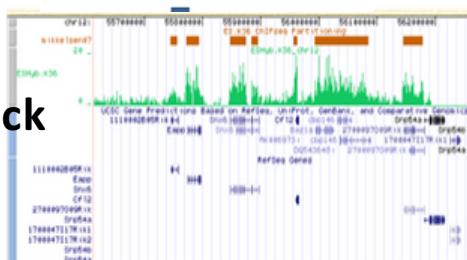
Quality control
fastQC



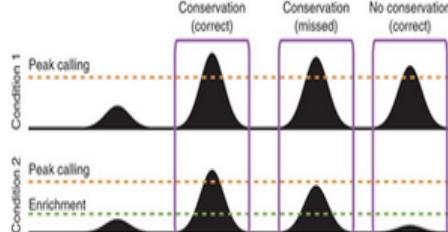
Reads mapping
bowtie



Build coverage track
perl



Peak calling
GRange, MACS



Raw data (sequence with quality scores)

Check quality of reads

Map and align reads to genome

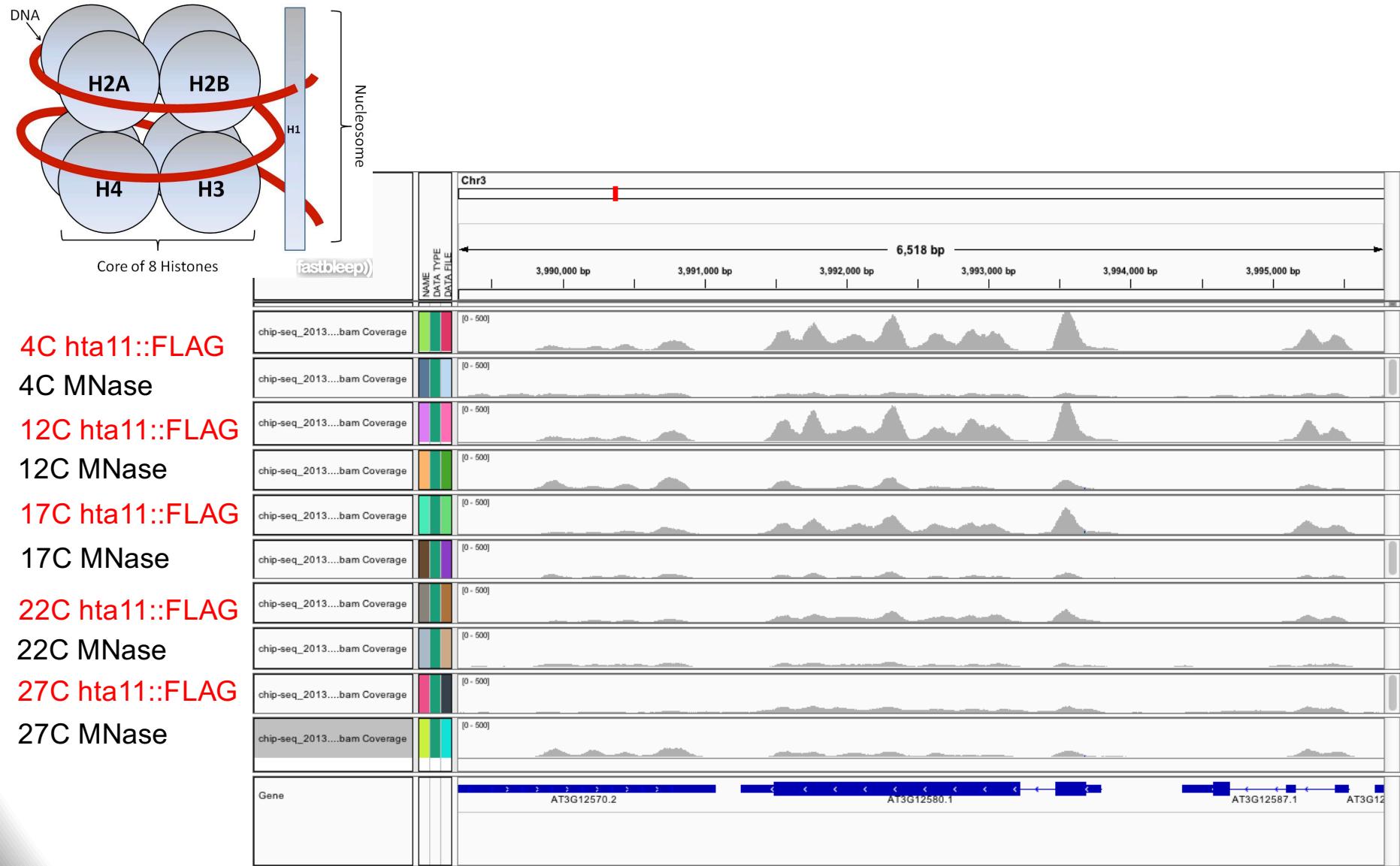
Visualize mapped data

Statistical analysis

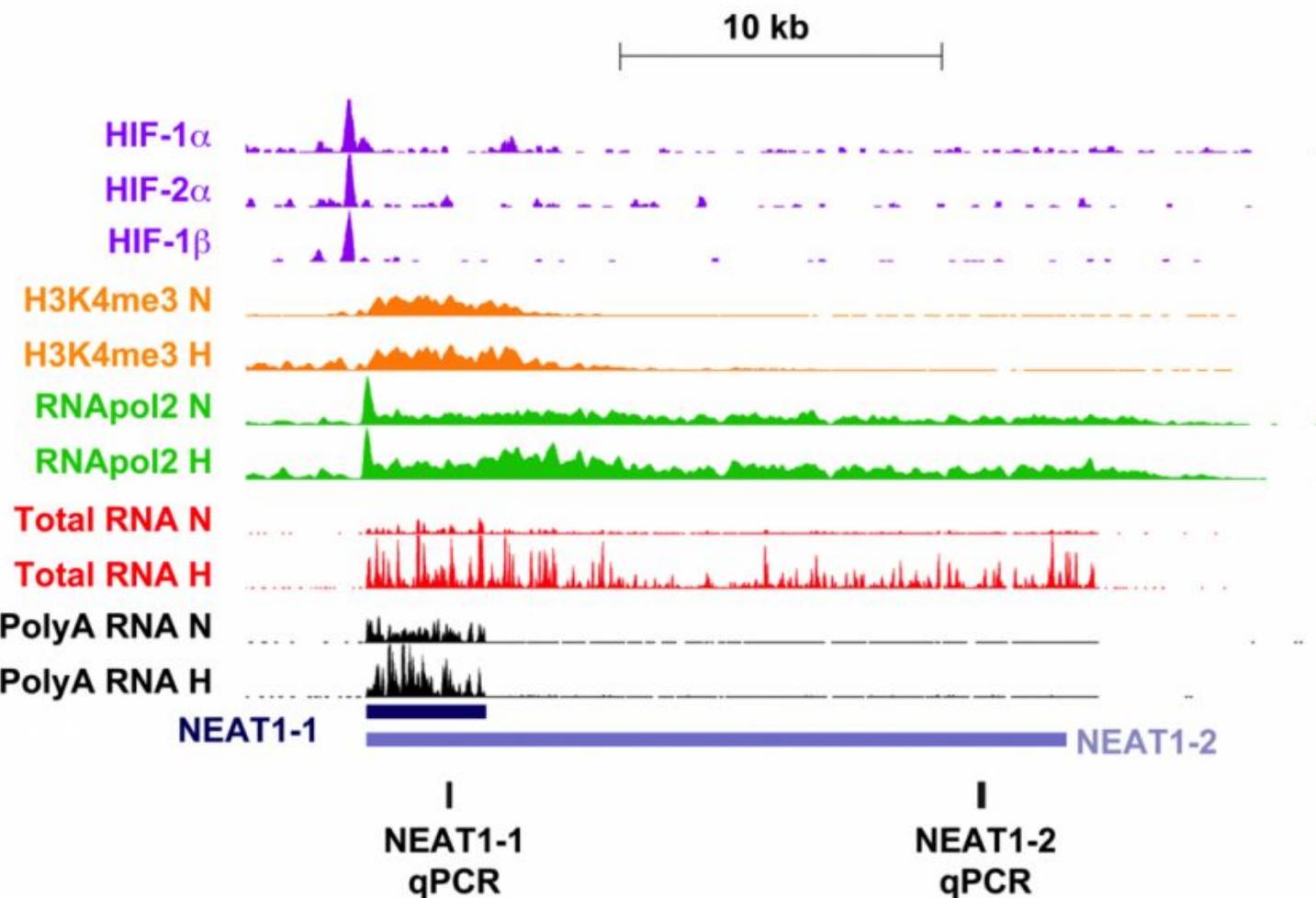
Adapted from

http://www4.utsouthwestern.edu/mcdermottlab/NGS/analysis/analysis_chipseq.html

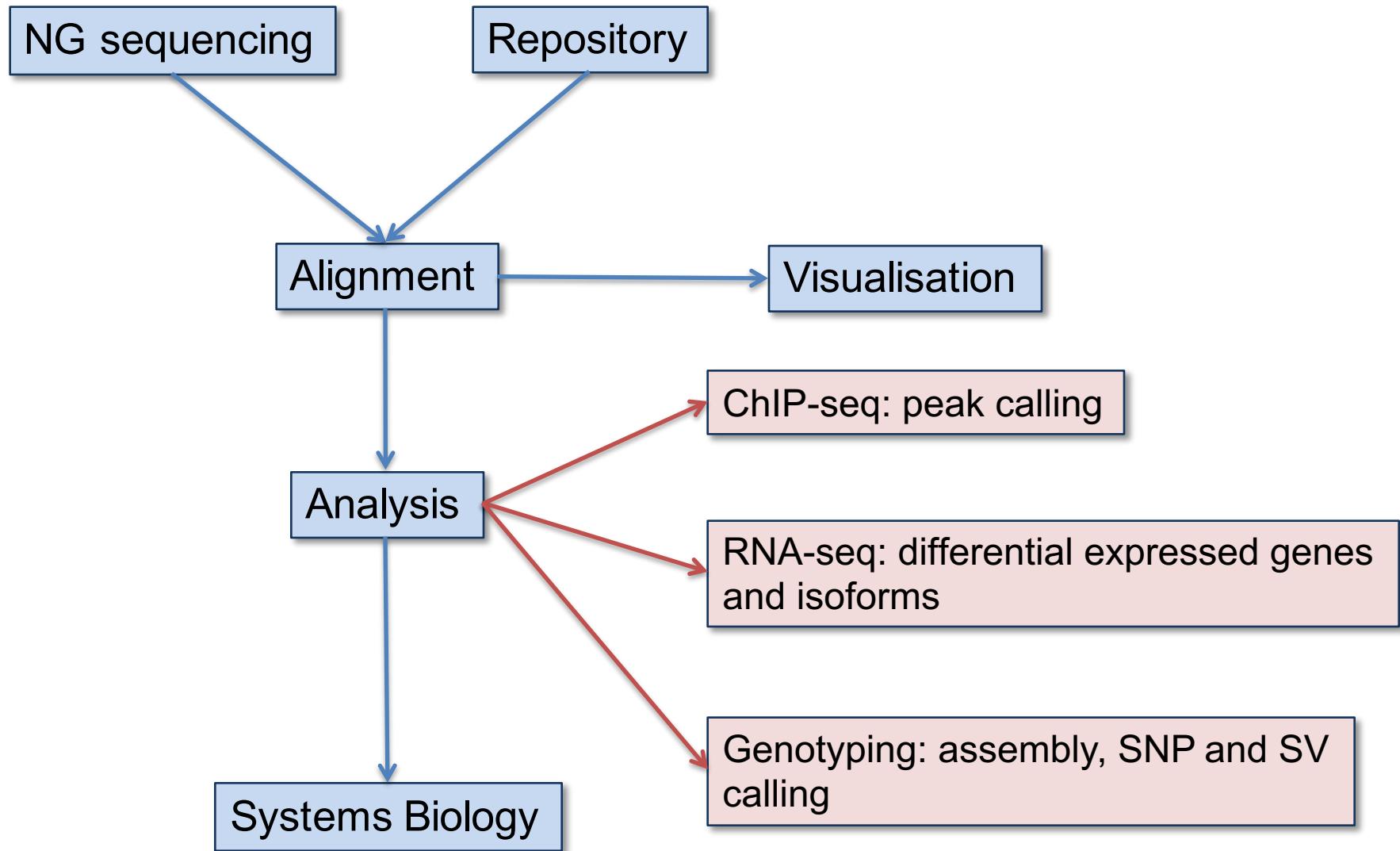
Examples of ChIP-data: H2A.Z position at HSP70 promoter



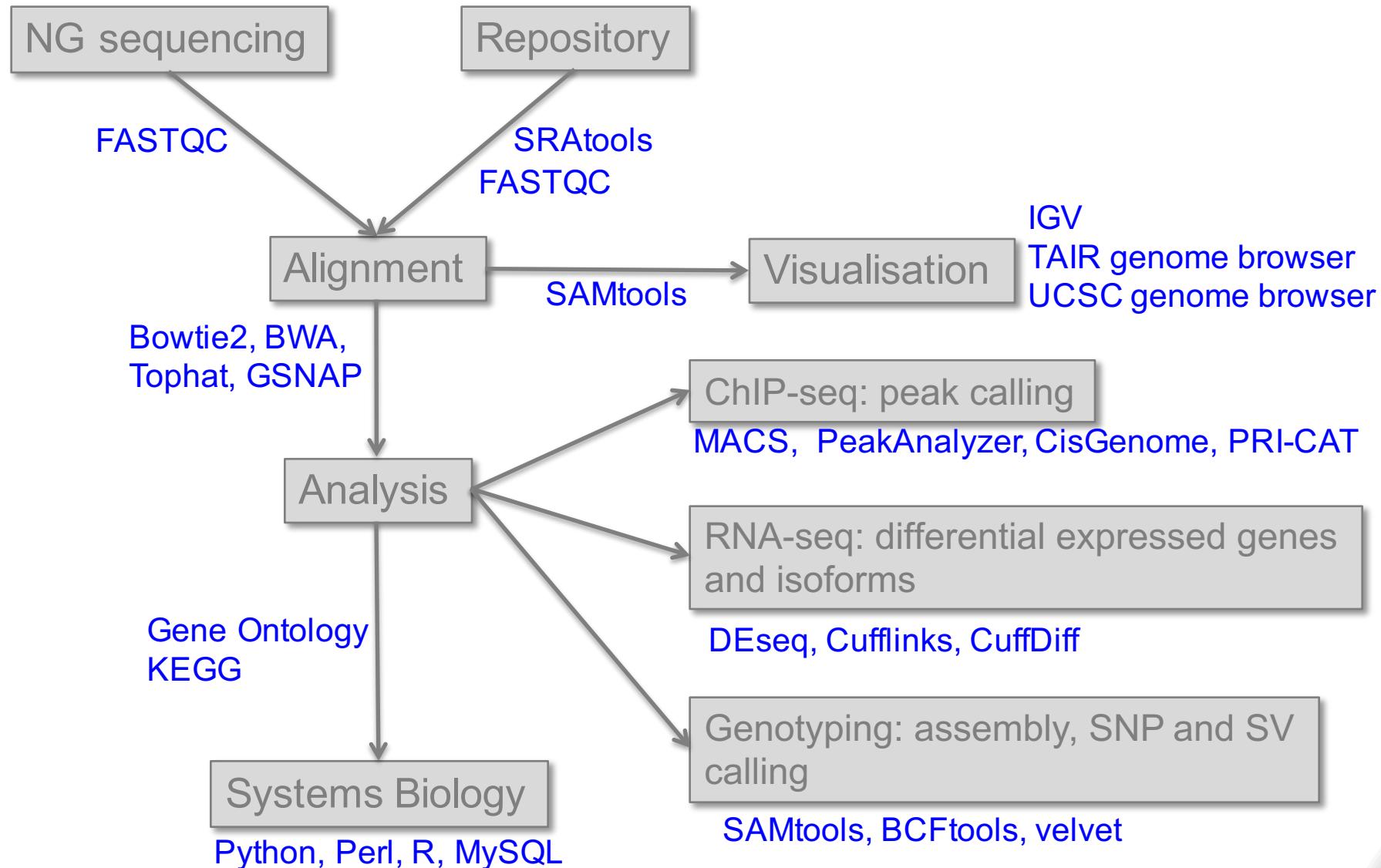
Examples of ChIP-data: transcriptional and epigenetic response to hypoxia in normal and cancer cells



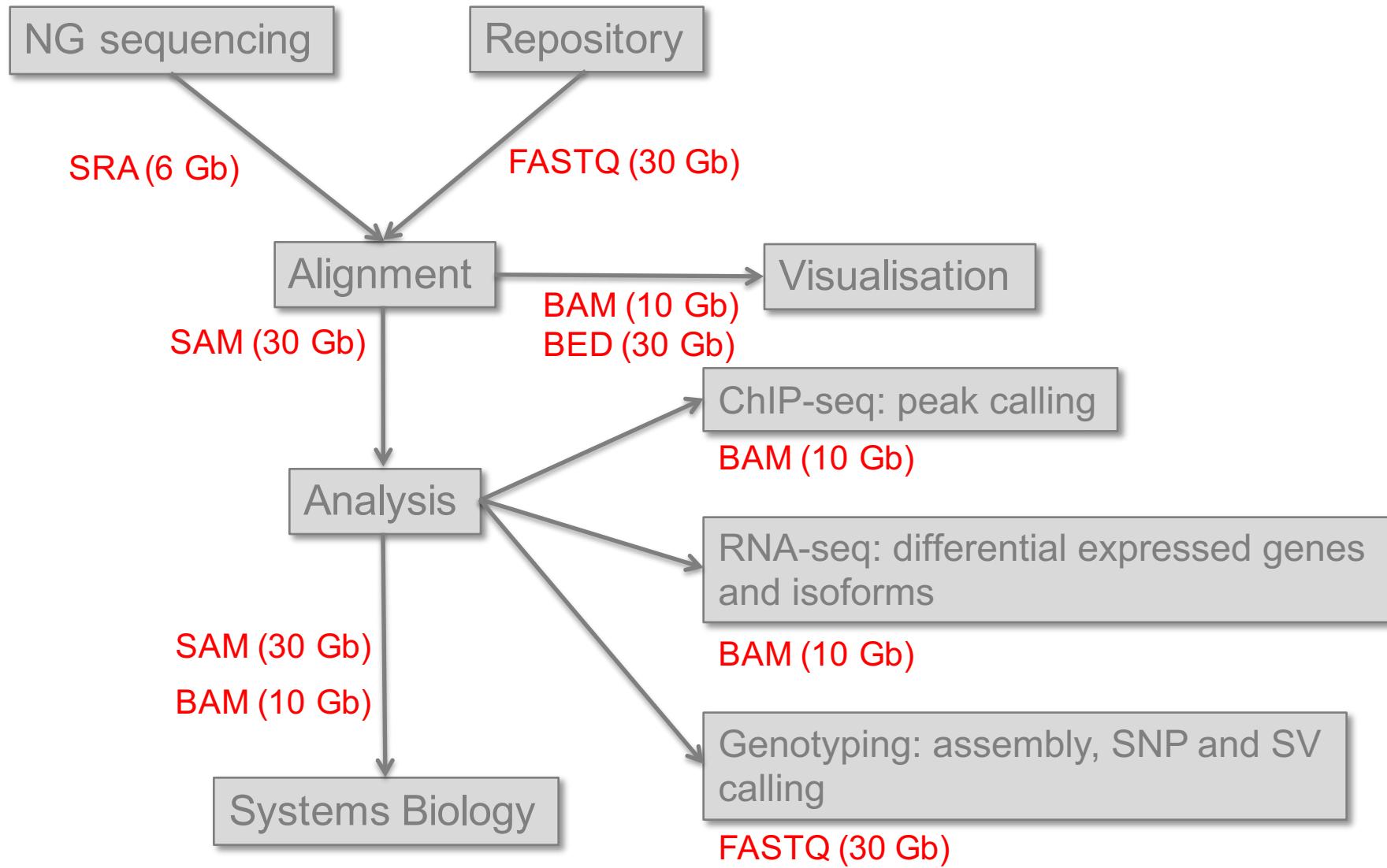
Computational Steps of NGS Data Analyses



Computational Steps of NGS Data Analyses



Computational Steps of NGS Data Analyses



FASTQ and Phred score

```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTT
+
! ' ' * ( ( ( ***+ ) ) % % + + ) ( % % % ) . 1 * * * - + * ' ' ) ) **55CCF>>>>CCCCCCC65
```

$$Q = -10 \log_{10} P$$

or

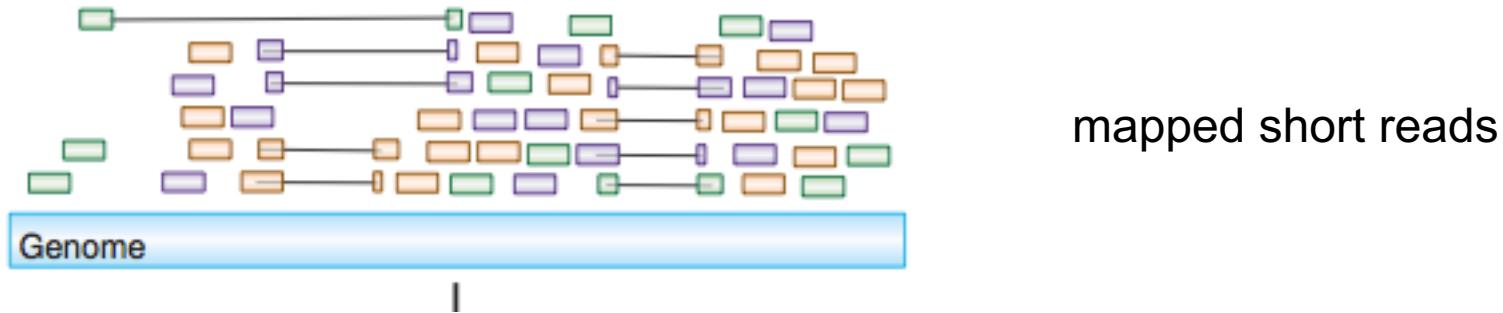
$$P = 10^{\frac{-Q}{10}}$$

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10000	99.99%
50	1 in 100000	99.999%

Phred score + 33 -> ASCII character

SAM format

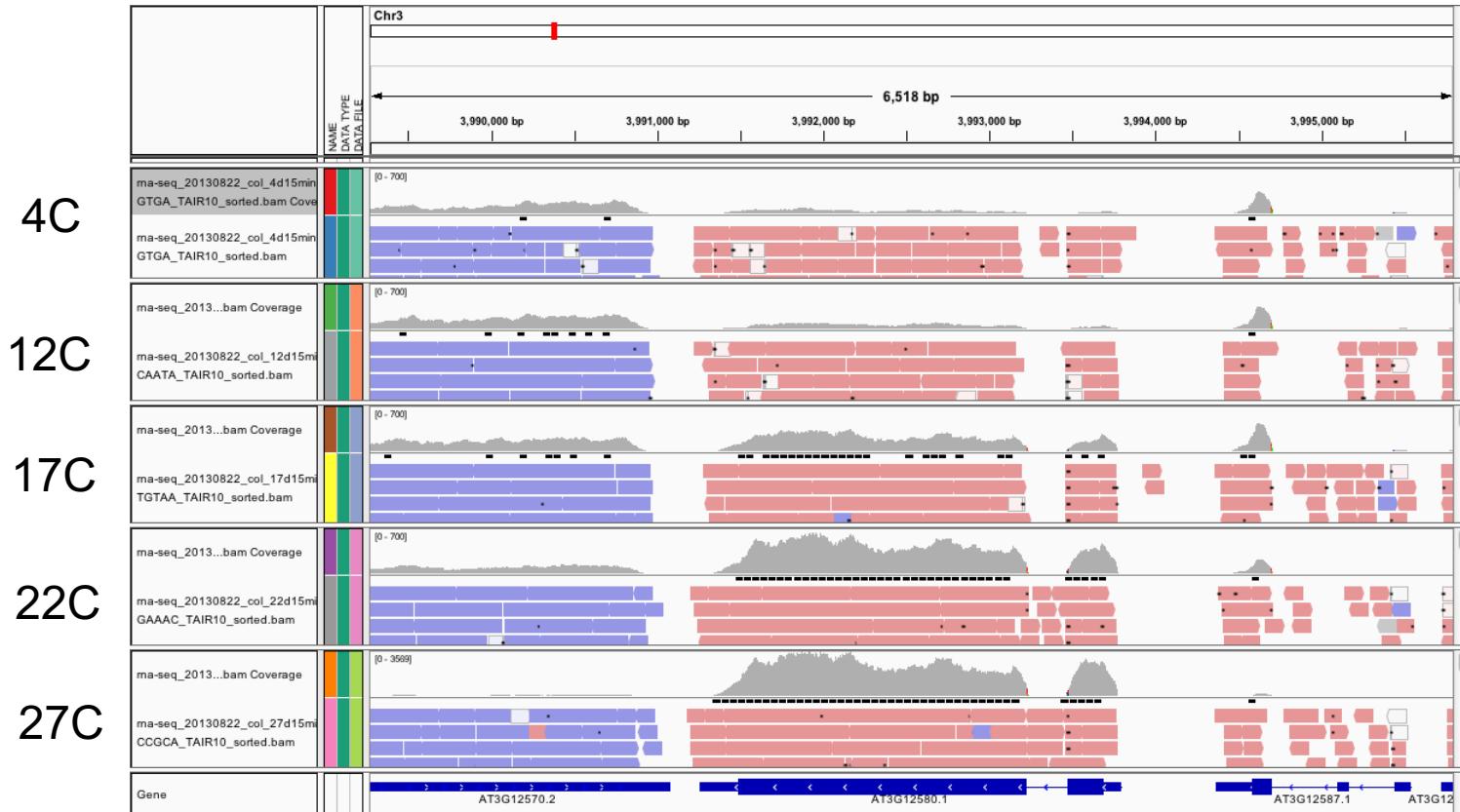


```
1 : 497 : R : -272+13M17D24M 113      1    497      37  37M 15
100338662 0
CGGGTCTGACCTGAGGAGAACTGTGCTCCGCCTTCAG
0 ;=====9 ;>>>>=>>>>>>>=>>>>>>>>
XT:A:U NM:i:0 SM:i:37      AM:i:0 X0:i:1 X1:i:0 XM:i:0
XO:i:0 XG:i:0 MD:Z:37
```

Outlines

- Transcriptomics: genome-wide transcriptional read-outs (**microarray, RNA-seq**)
- Regulation of transcription (and transcriptomes)
 - Computational prediction
 - *In vitro*: EMSA, PBM
 - *In vivo*: ChIP-chip, ChIP-seq
- Putting things in perspectives

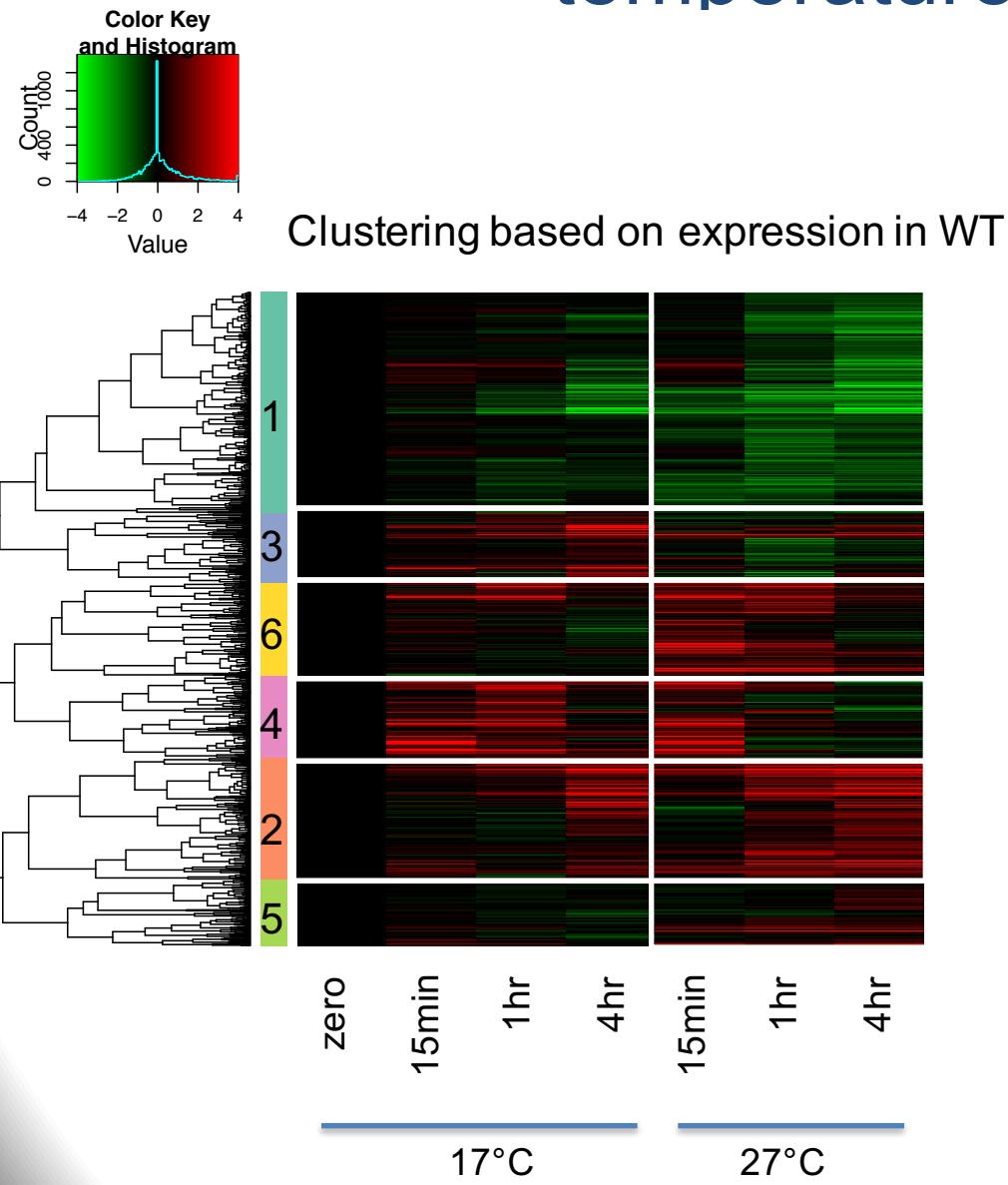
Examples of RNA-seq data: HSP70 as a temperature sensor



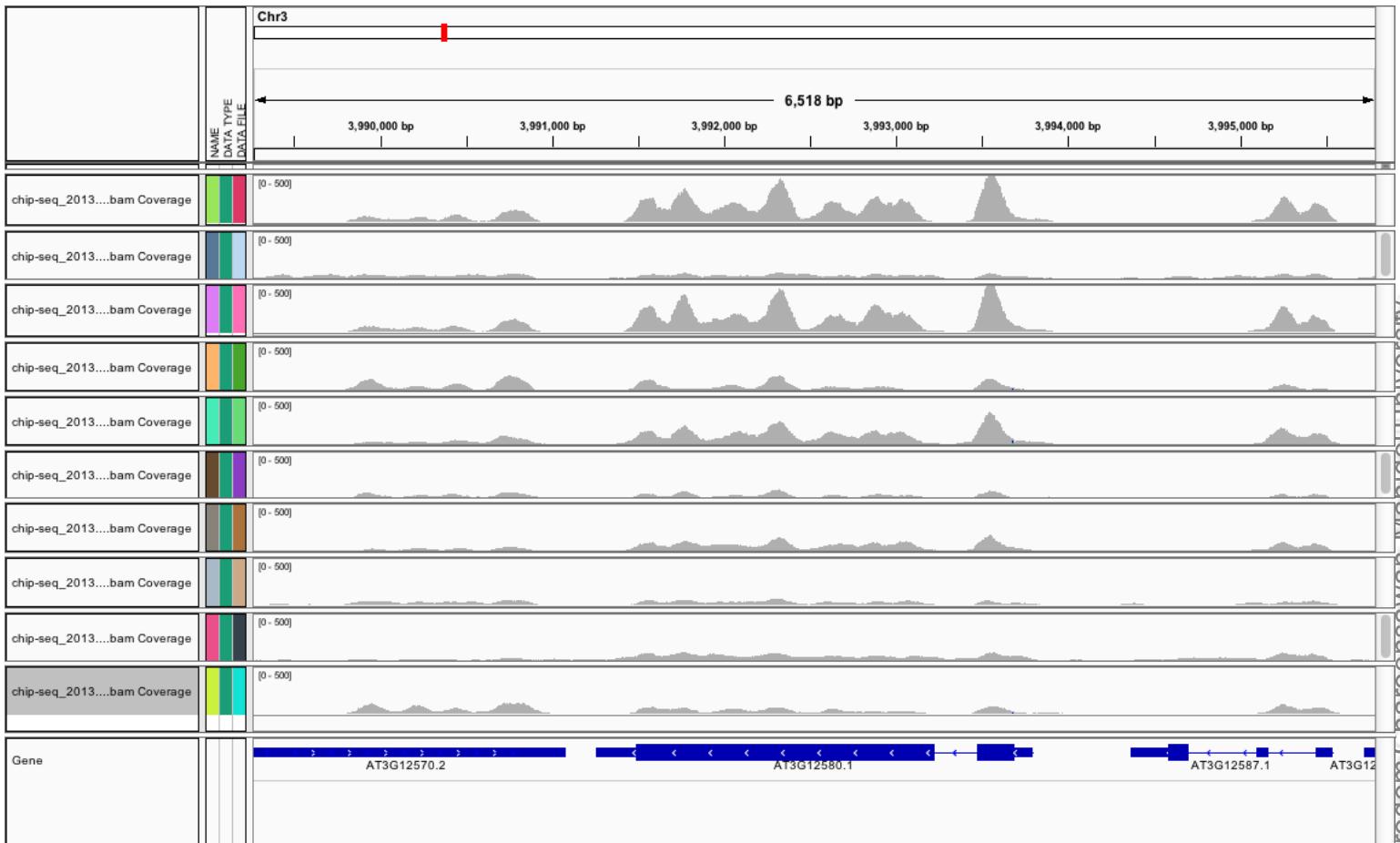
Transcriptional response after 15 mins

With Sandra Cortijo and Phil Wigge (Scale 0-600 for all, except 27C: 0-3500)

Transcriptome dynamics in response to temperature variation

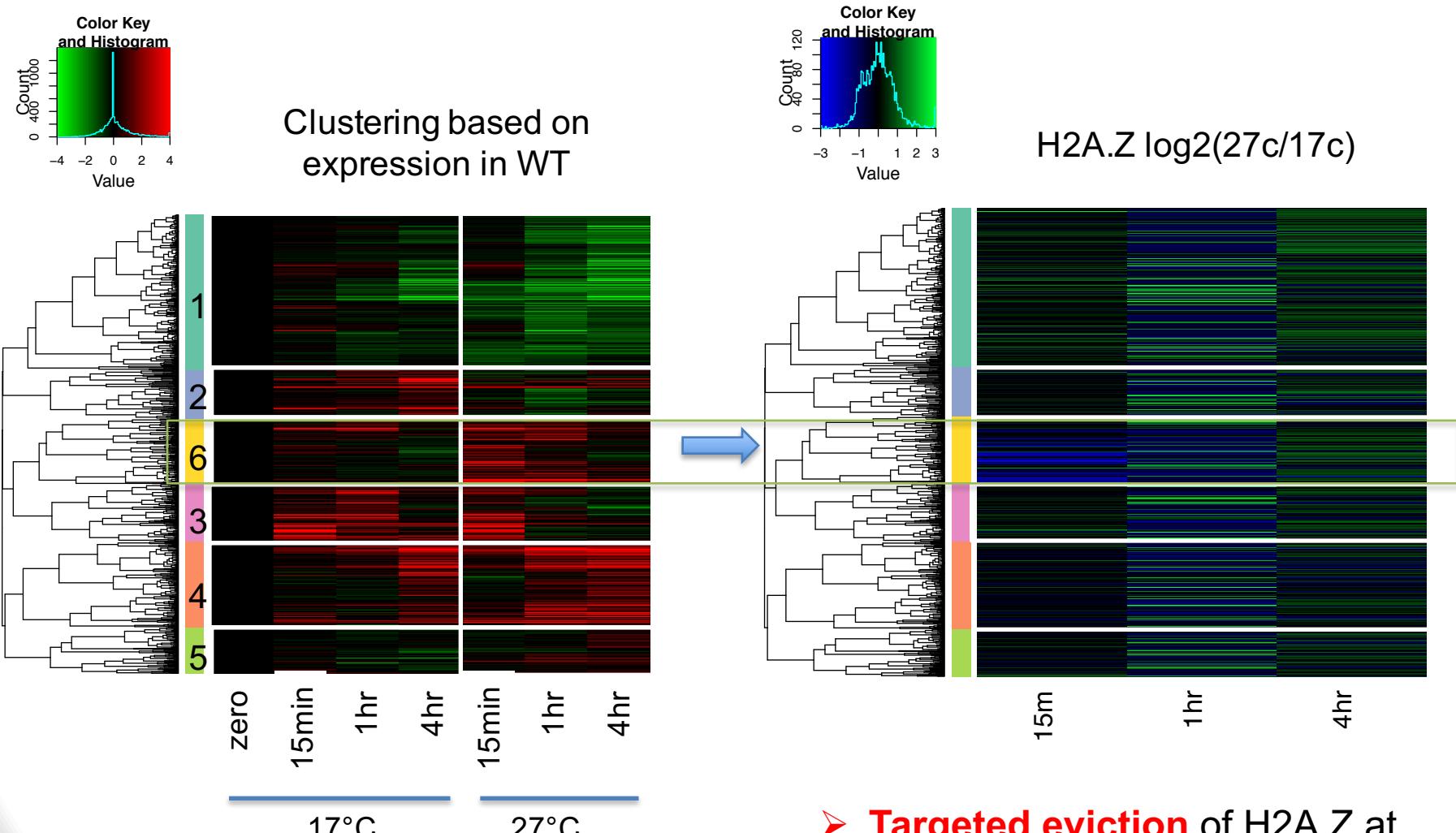


Examples of ChIP-data: H2A.Z position at HSP70 promoter



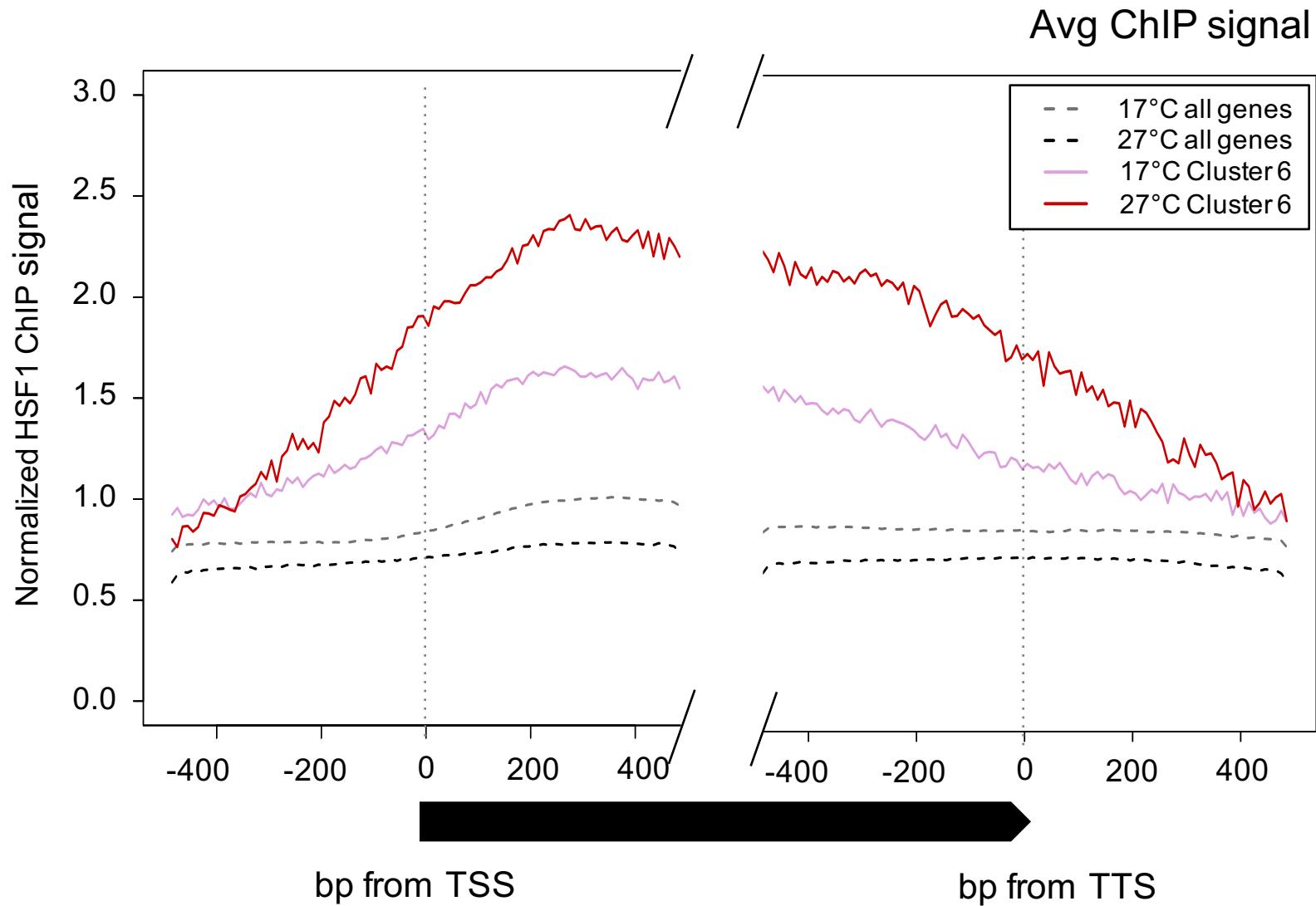
H2AZ occupancy after 1 hr Scales 0-500 for all

Interplay between H2A.Z and transcriptional changes in response to temperature



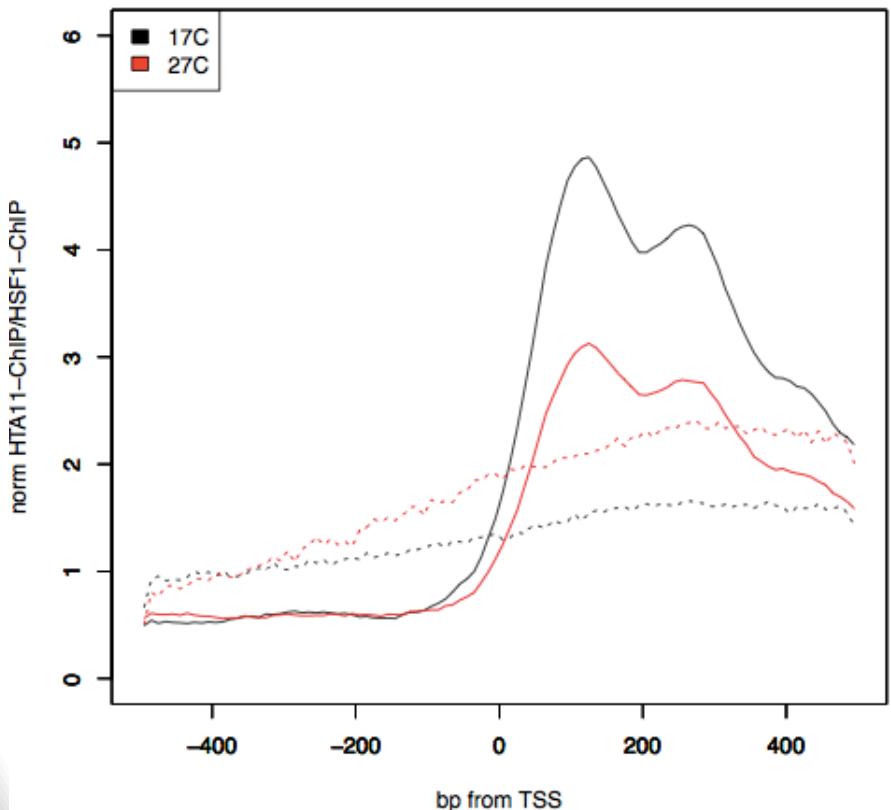
➤ **Targeted eviction** of H2A.Z at 15min at genes activated by the temperature shift

Increase of HSF1 binding at highly expressed genes in Cluster 6

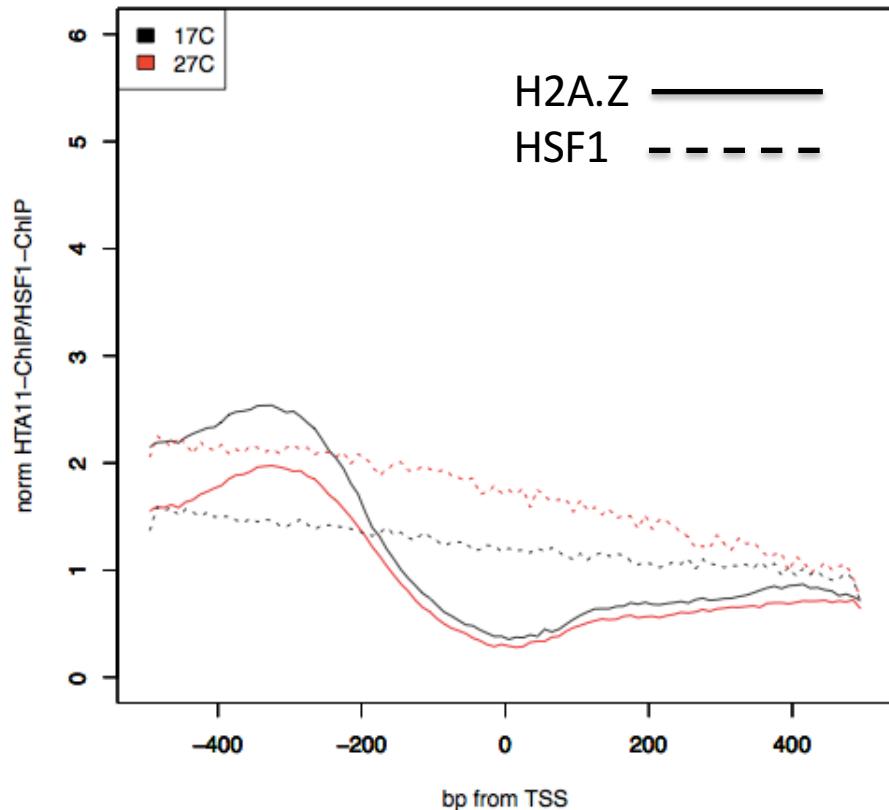


Reduction of nucleosome stability vs' increase of HSF1 binding at highly expressed genes in Cluster 6

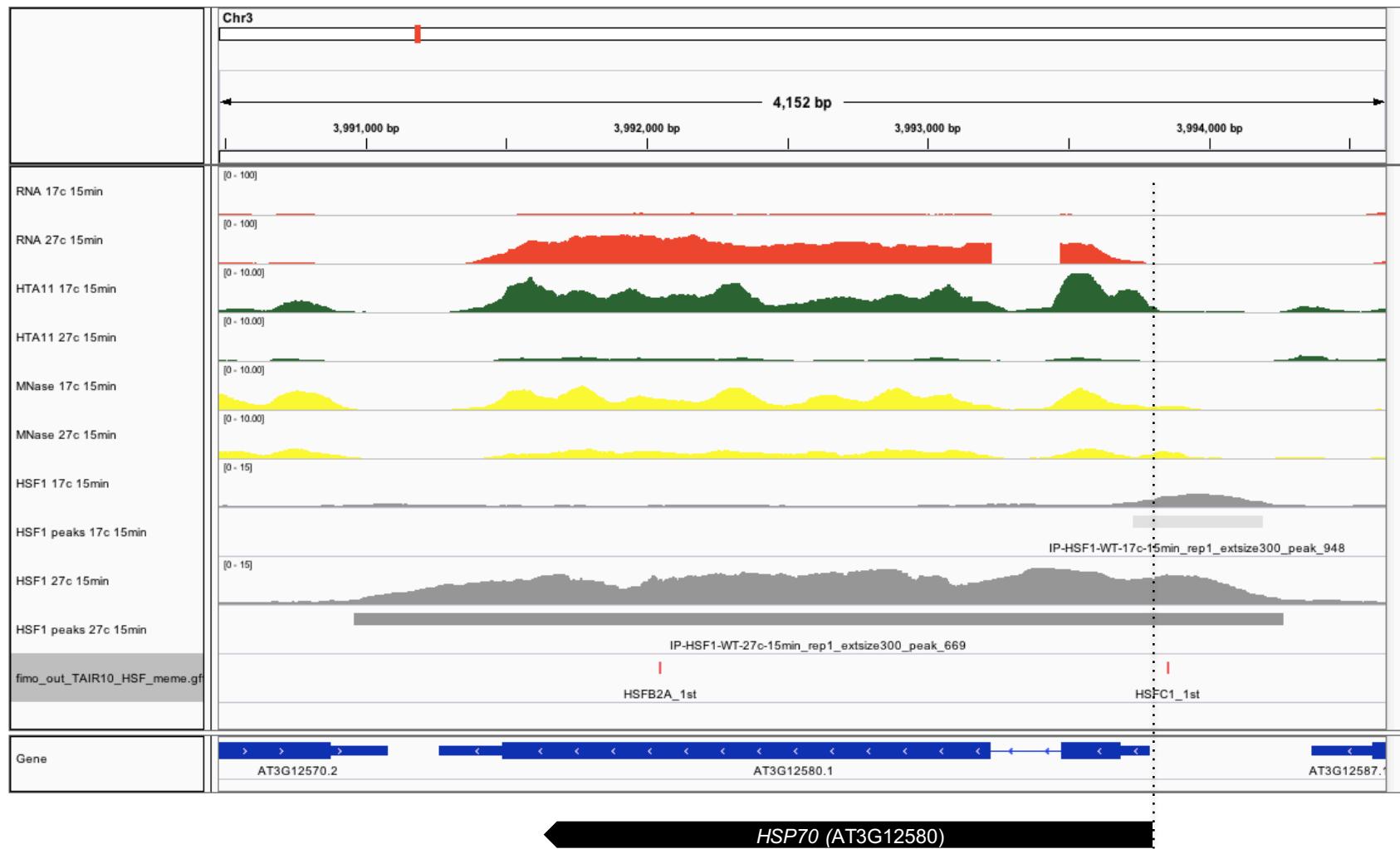
Cluster6 protein coding genes – 15min



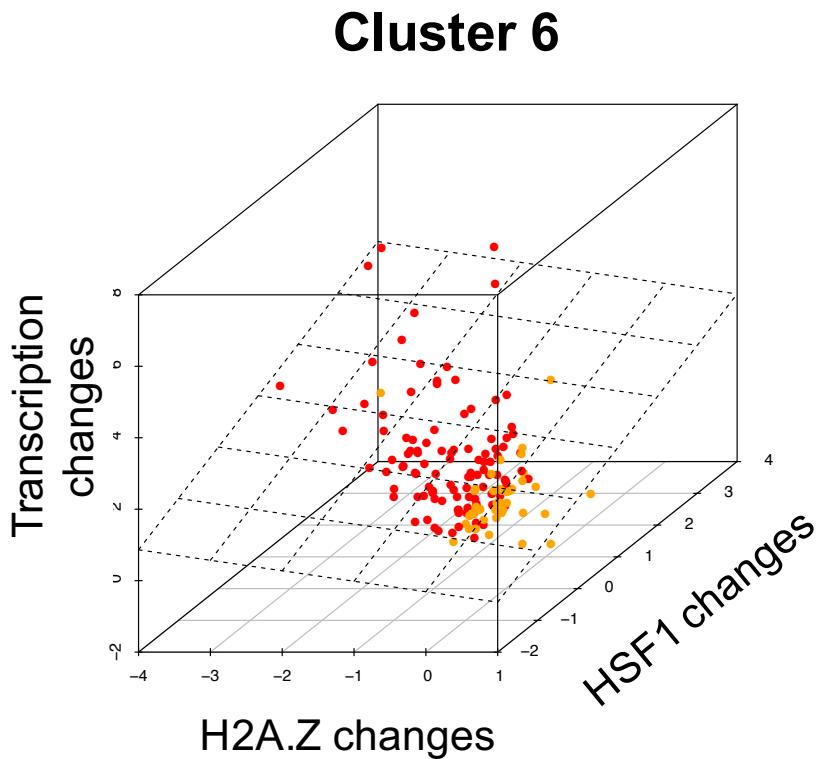
Cluster6 protein coding genes – 15min



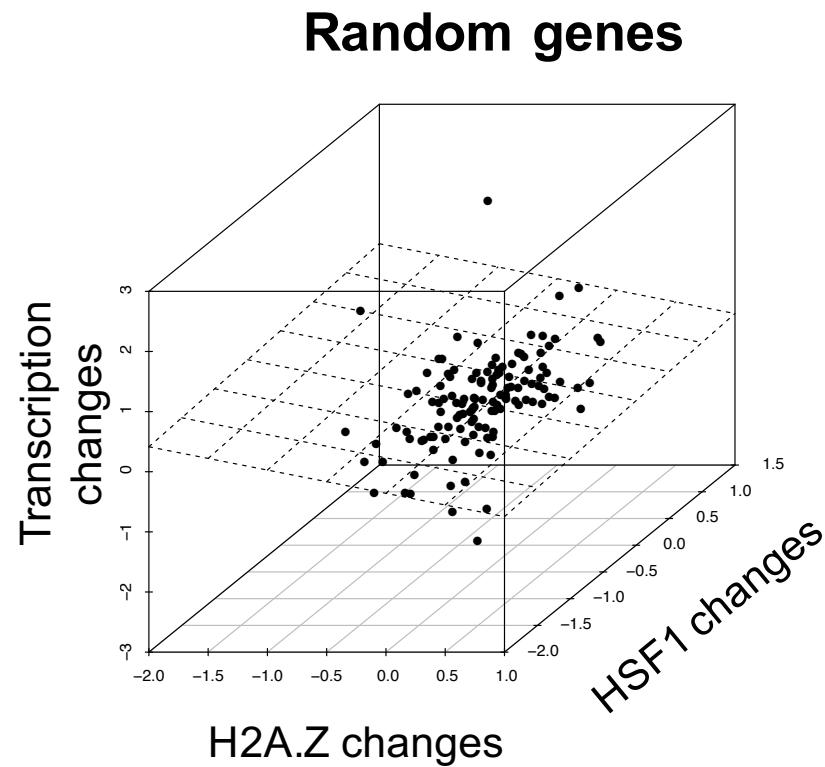
Antagonistic role of HSF1 and H2A.Z



HSF1 binding and H2A.Z eviction are good predictors for transcription of temperature-responsive genes



Multiple R-squared: 0.33
P-value: 3.678e-13



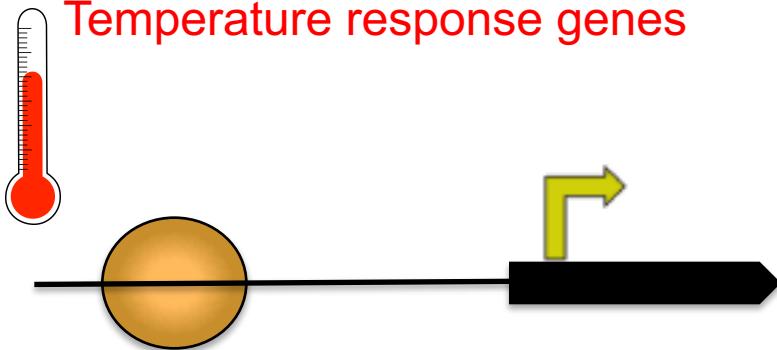
Multiple R-squared: 0.06
P-value: 0.0221

Characterising Transcriptional Architectures of Temperature Response



Thale cress

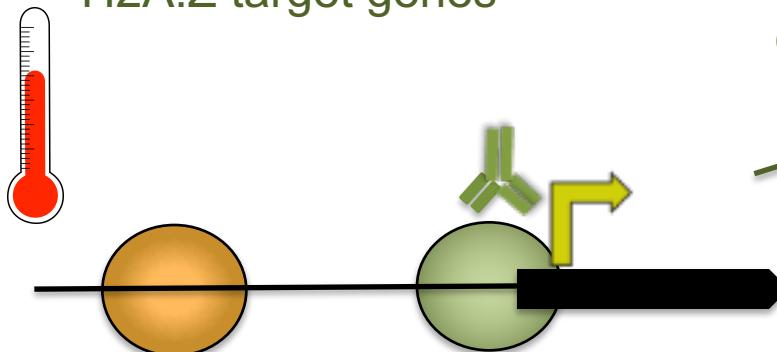
Temperature response genes



RNA-seq



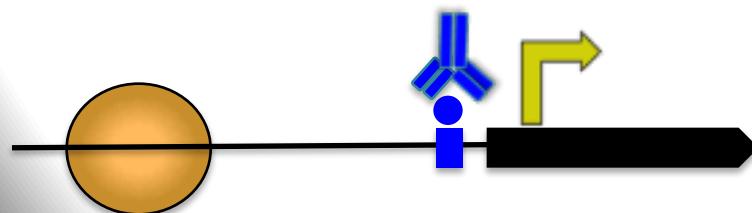
H2A.Z target genes



ChIP-seq



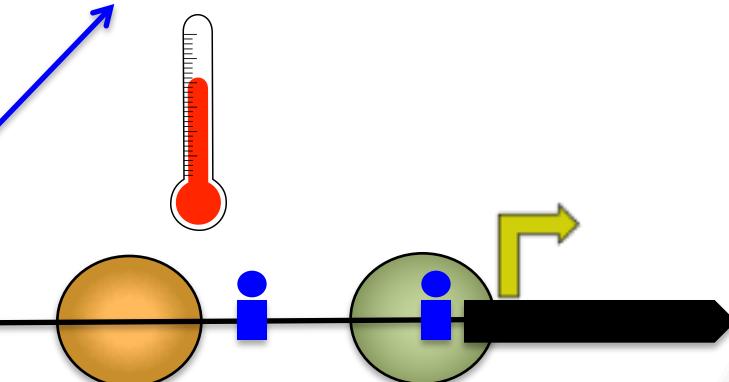
Trans-acting factors' (TFs) target genes



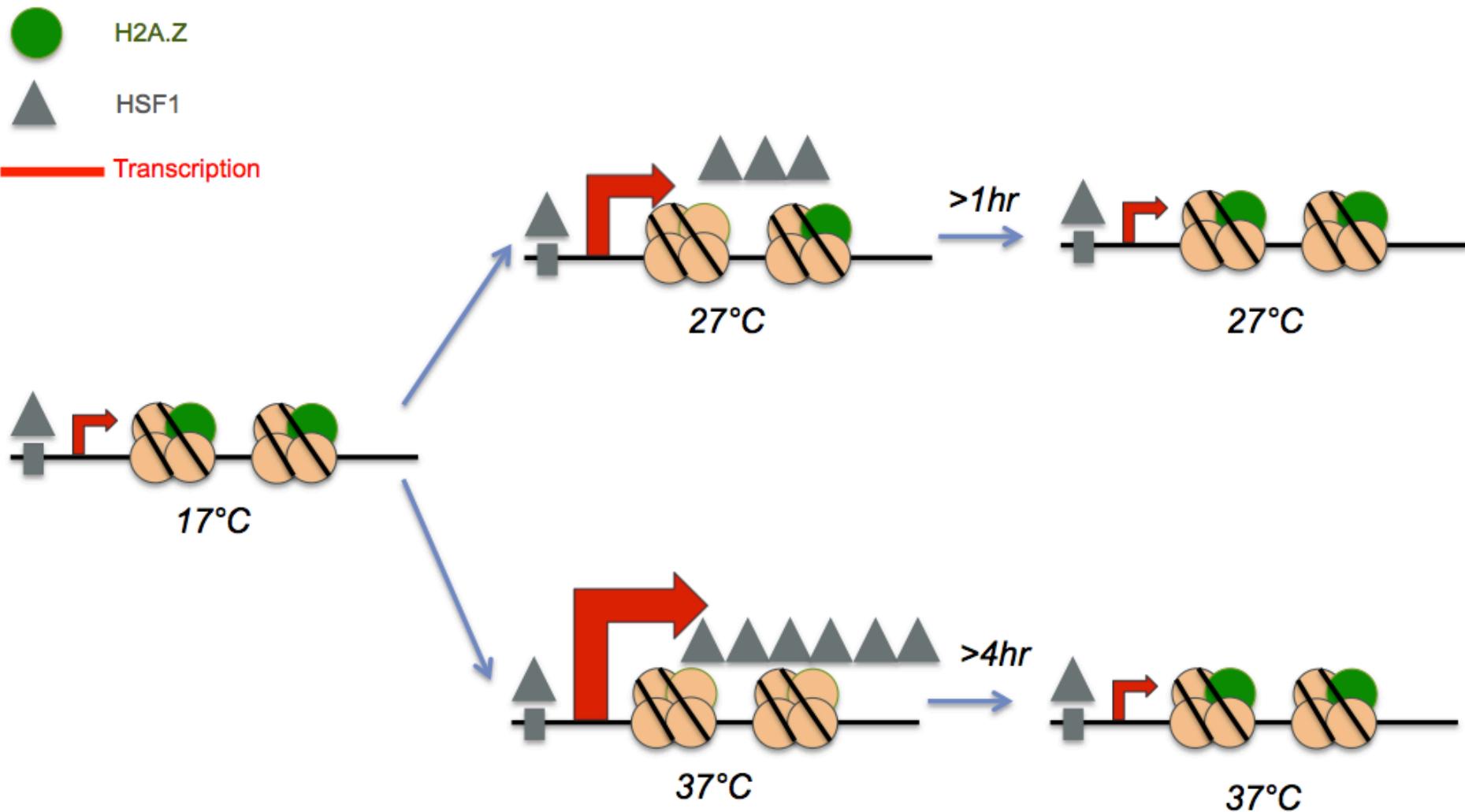
ChIP-seq



How many? What genes?

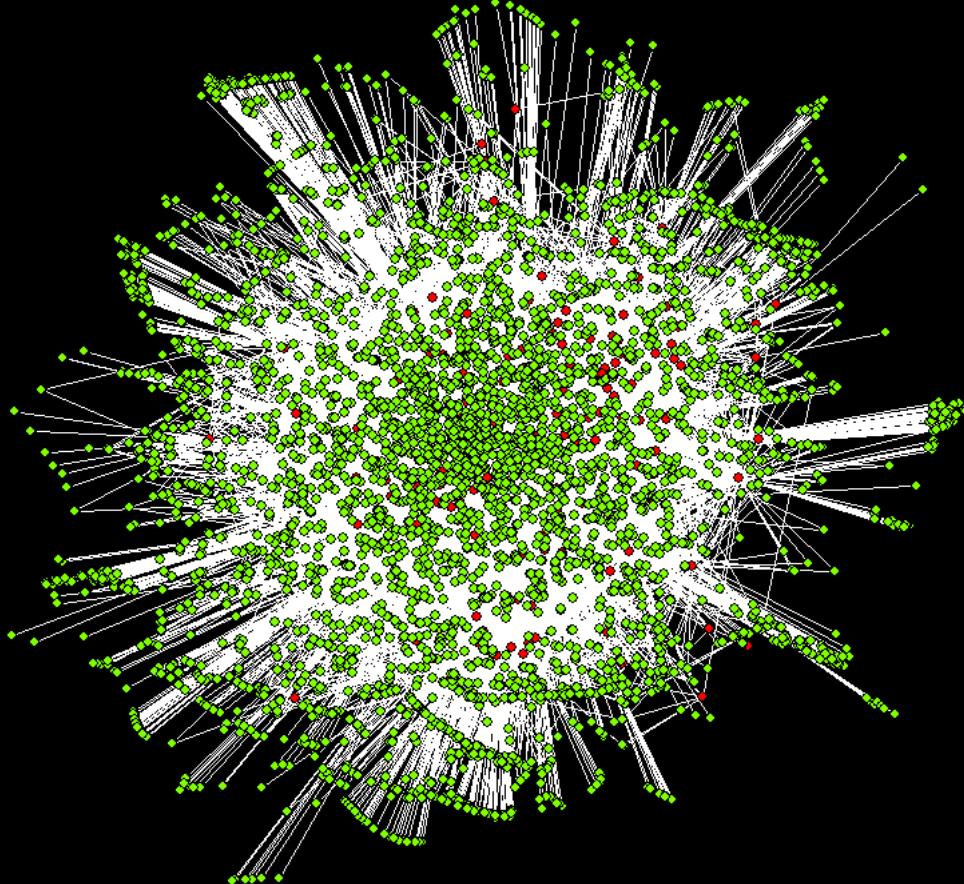


How temperature affects gene expression and plant development?



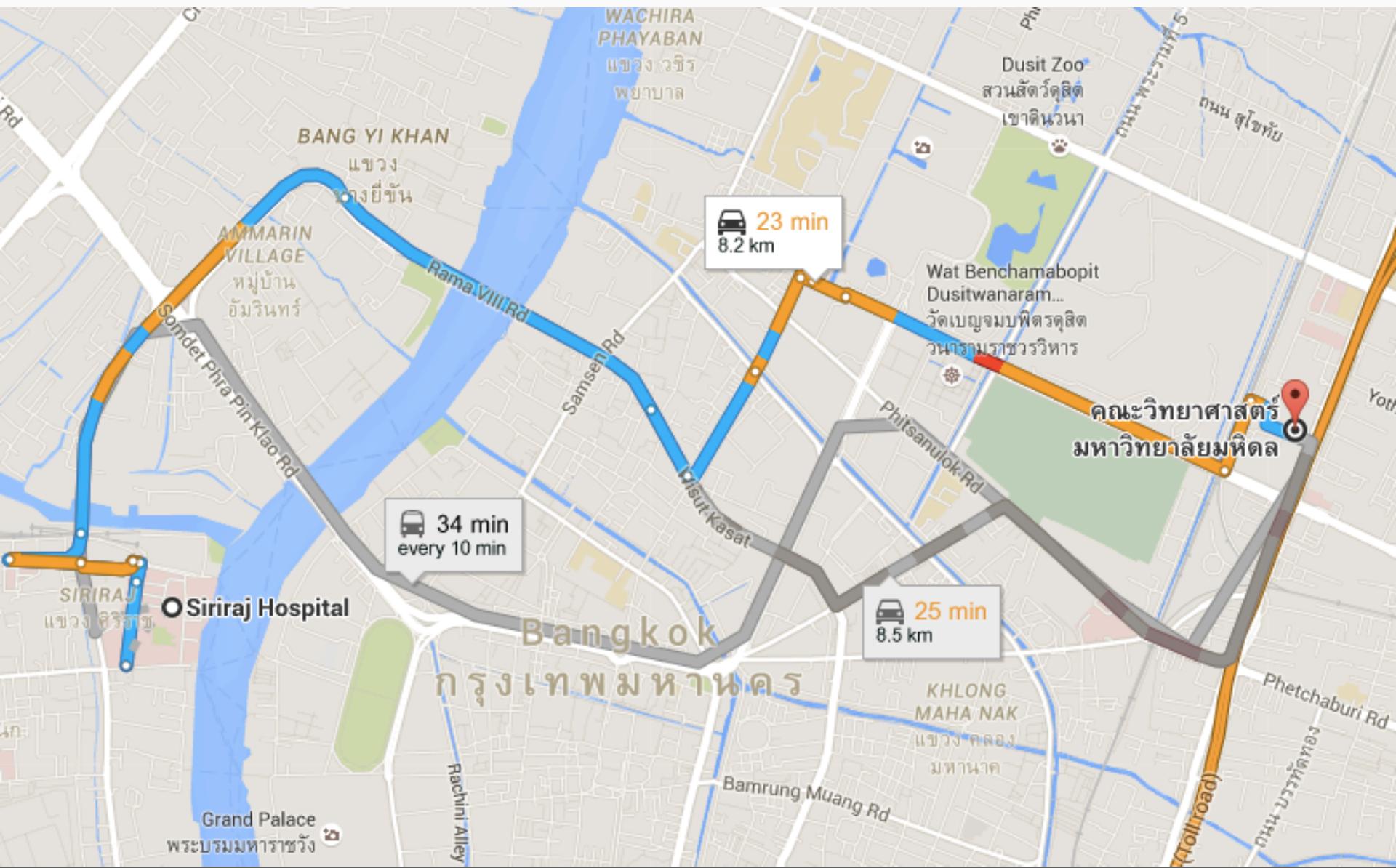
Cortijo* and Charoensawan* et al. Molecular Plant. 2017

Network Biology



Slides courtesy of Dr MM Babu, MRC LMB

How to from A to B?



Further readings

- Wang Z, Gerstein M, Snyder M. [RNA-Seq: a revolutionary tool for transcriptomics](#). *Nat Rev Genet*. 2009 Jan;10(1):57-63
- Furey TS. [ChIP-seq and beyond: new and improved methodologies to detect and characterize protein-DNA interactions](#). *Nat Rev Genet*. 2012 Dec;13(12):840-52.
- Pepke S, Wold B, Mortazavi A. [Computation for ChIP-seq and RNA-seq studies](#). *Nat Methods*. 2009 Nov;6(11 Suppl):S22-32. doi: 10.1038/nmeth.1371.
- Hrdlickova R et al. [RNA-Seq methods for transcriptome analysis](#). *Wiley Interdiscip Rev RNA*. 2017 Jan;8(1). doi: 10.1002/wrna.1364

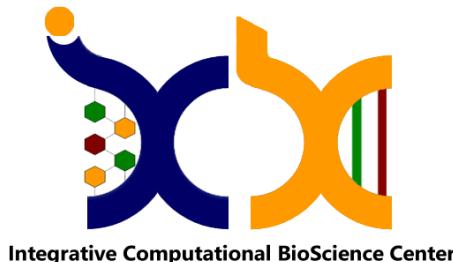
Question? Comment?

Varodom Charoensawan

Department of Biochemistry

Faculty of Science, Mahidol University

varodomc@gmail.com, varodom.cha@mahidol.ac.th



Mahidol
University
Wisdom of the Land